

UNIVERSITÄT
BAYREUTH



Foundation Models

Generative AI

Simeon Allmendinger

University of Bayreuth
Project ABBA

Chair for Information Systems and Human-Centric AI

Fraunhofer Institute for Applied Information Technology FIT

www.uni-bayreuth.de | www.fit.fraunhofer.de

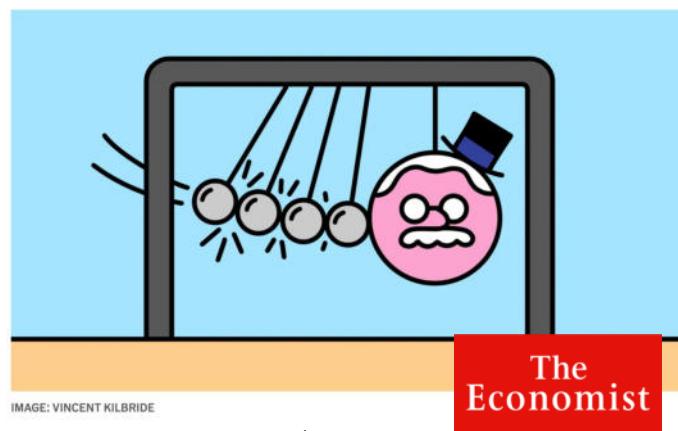
Is there a concentration of power in AI?

Foundation Models

Leaders | Undisrupted

AI could fortify big business, not upend it

Upstarts face an uphill battle



[1]



European Commission - Speech
[Check Against Delivery]



[2]

MAKING ARTIFICIAL INTELLIGENCE AVAILABLE TO ALL - HOW TO AVOID BIG TECH'S MONOPOLY ON AI?

Brussels, 19 February 2024

- Check against delivery -

Good afternoon.

Let me start by thanking Stéphanie Yon-Courtin for organising this important event. As Competition Commissioner, I have been very pleased with your interest and dedication to competition policies over the last five years in the European Parliament.

Thank you for putting the spotlight on Artificial Intelligence and competition. Because '*human* intelligence' is exactly what we need right now, to strike the right balance on intelligence of the *artificial* kind. In order to shape the emerging markets that are enabled by Large Language Models and other applications in AI. To make sure that competition can thrive, and consumers reap the benefits of these new markets, without hampering their development.

By thinking ahead, by acting swiftly and by cooperating, we have a window of opportunity to maximise these benefits; while at the same time, minimising the risks. But that window is closing. If we don't act soon, we will find ourselves, once again, chasing solutions to problems we did not anticipate. So debates like this one are not only very timely, they are also urgent.

February 19th, 2024

<https://www.economist.com/leaders/2023/08/24/ai-could-fortify-big-business-not-upend-it> [1]
https://ec.europa.eu/newsroom/ecpc-failover/pdf/speech-74-931_en.pdf [2]

Hello!



Simeon Allmendinger
PhD Student

**Chair for Information
Systems and human-centric
Artificial Intelligence**



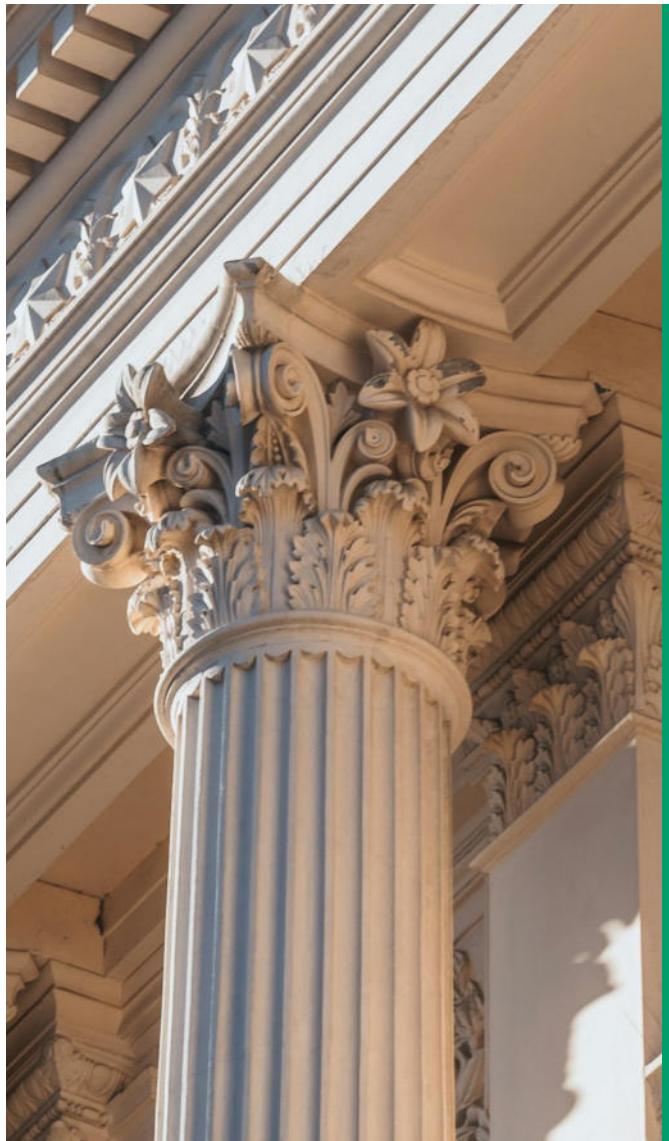
**Production /
Manufacturing**



**Autonomous
Driving**



**Healthcare /
Surgery**



- 1 Introduction
- 2 Multimodality in Data
- 3 Generative Foundation Models
- 4 Research
- 5 Industry

What is the idea of Foundation Models?

A paradigm shift in machine learning enabled by transfer learning and scale.



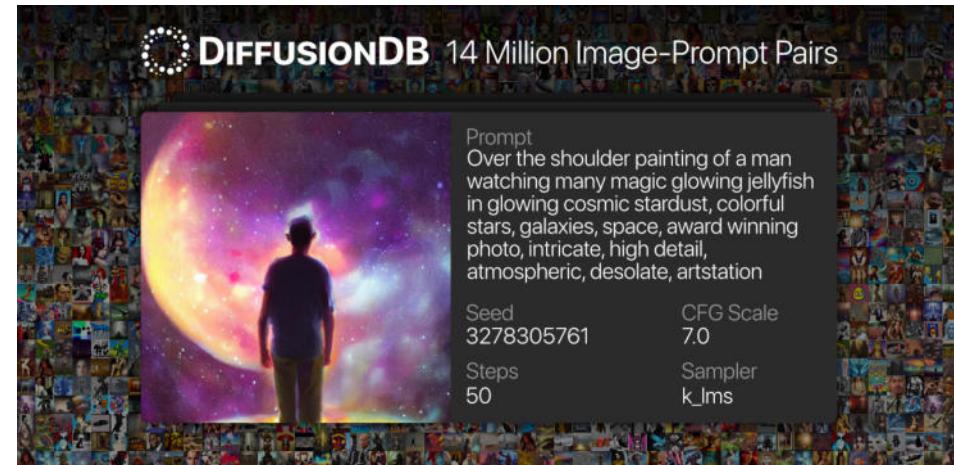
What are Foundation Models?

Train one model on a huge amount of data and adapt it to many applications.

Foundation Model



Huge datasets



[3]

<https://github.com/poloclub/diffusiondb> [3]

What are Foundation Models?

Train one model on a huge amount of data and adapt it to many applications.

Foundation Model



Huge datasets



Unlabeled data
self-supervised learning



the university has a beautiful campus

$y_{<t}$

History

y_t

Most probable
next token, e.g.,
word, prefix,
suffix, etc.

What are Foundation Models?

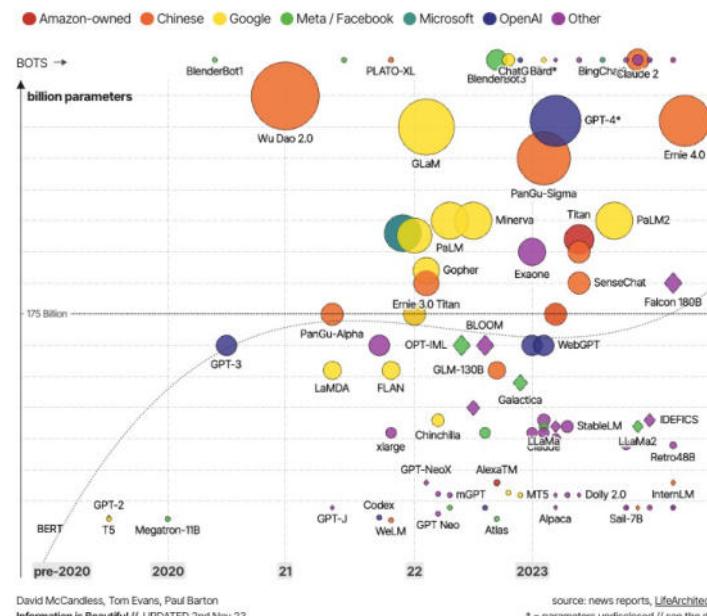
Train one model on a huge amount of data and adapt it to many applications.

Foundation Model

-  Huge datasets
-  Unlabeled data self-supervised learning
-  Large models



The Rise and Rise of A.I.
Large Language Models (LLMs) & their associated bots like ChatGPT



[4]

<https://informationisbeautiful.net/visualizations/the-rise-of-generative-ai-large-language-models-llms-like-chatgpt/> [4]

What are Foundation Models?

Train one model on a huge amount of data and adapt it to many applications.

Foundation Model

-  Huge datasets
-  Unlabeled data
self-supervised learning
-  Large models

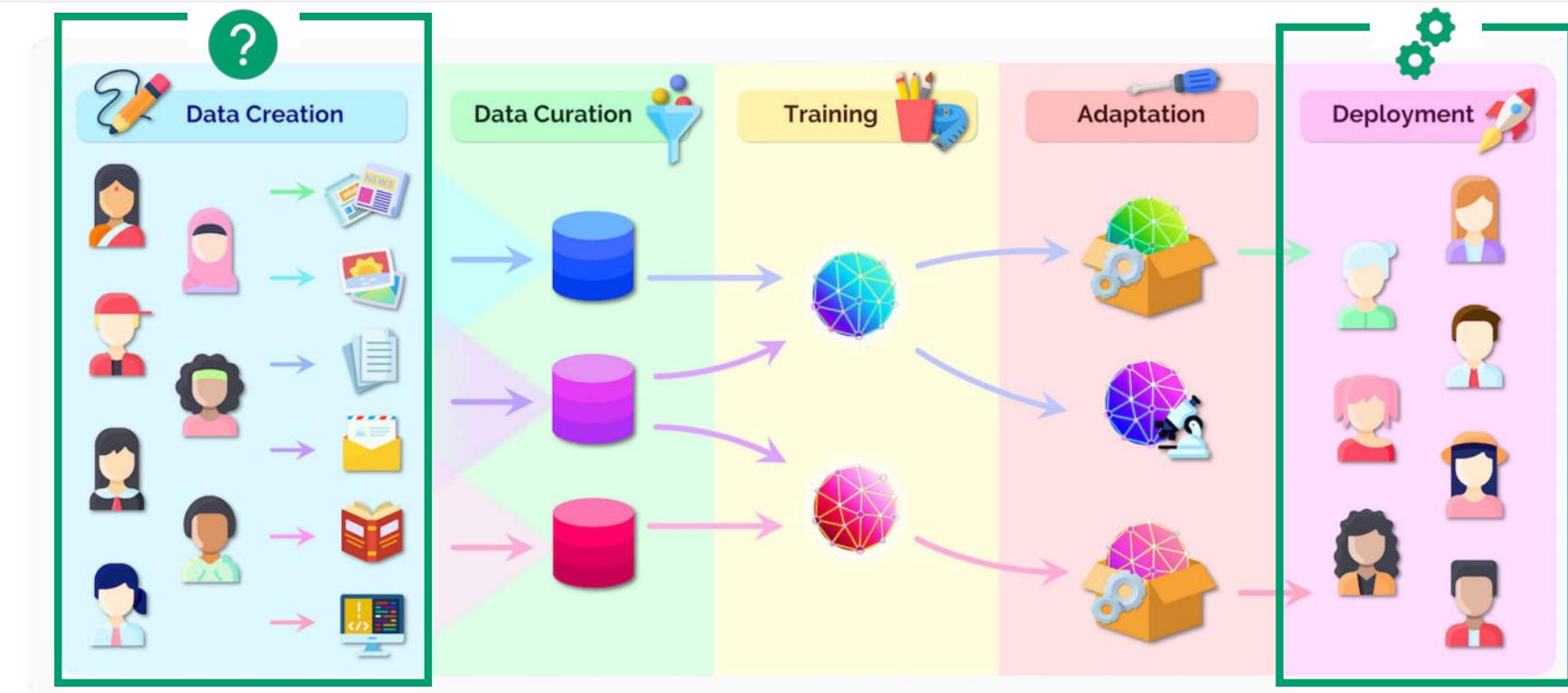
“

“a foundation model is itself incomplete but serves as the common basis from which many task-specific models are built via adaptation” [5]

Rishi Bommasani et al., (2022). On the Opportunities and Risks of Foundation Models. [5]

What are Foundation Models?

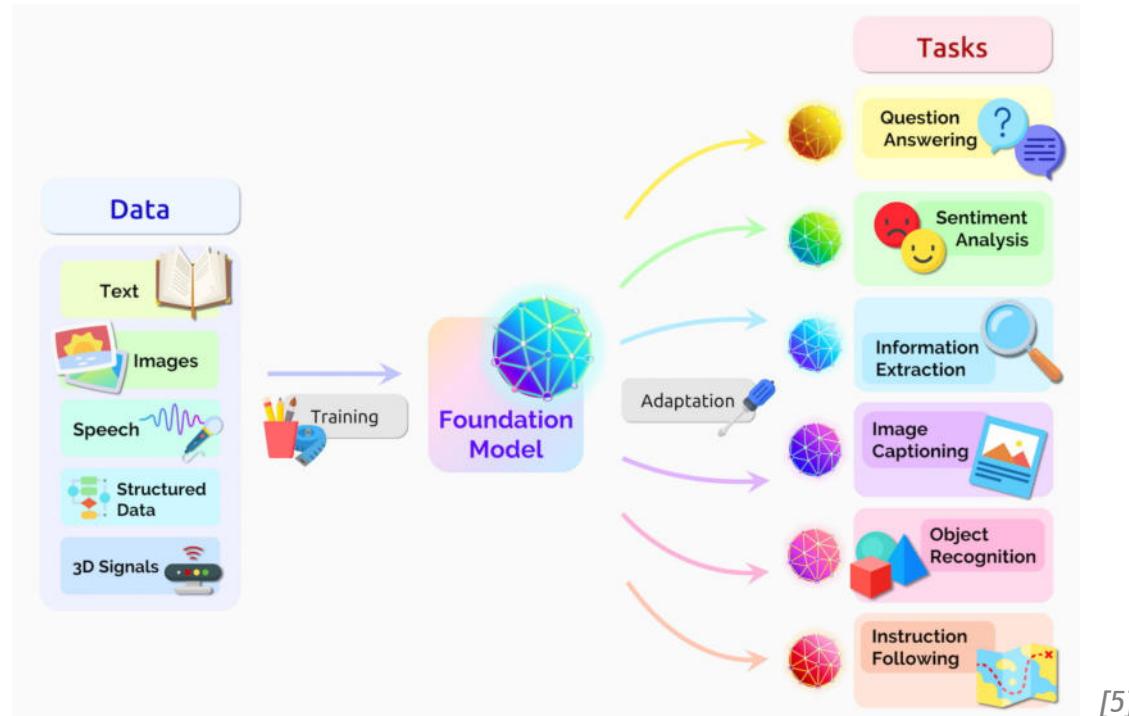
Train one model on a huge amount of data and adapt it to many applications.



Rishi Bommasani et al., (2022). On the Opportunities and Risks of Foundation Models. [5]

What are Foundation Models?

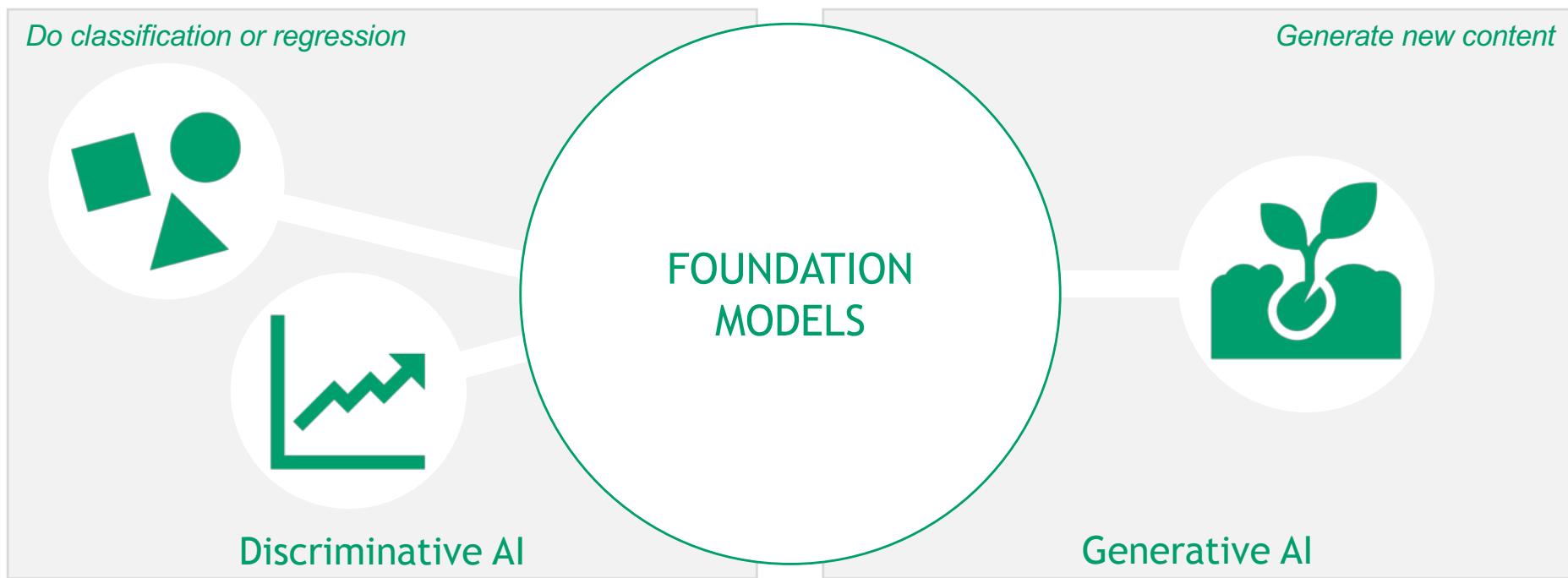
Foundation models go well beyond language.



Rishi Bommasani et al., (2022). On the Opportunities and Risks of Foundation Models. [5]

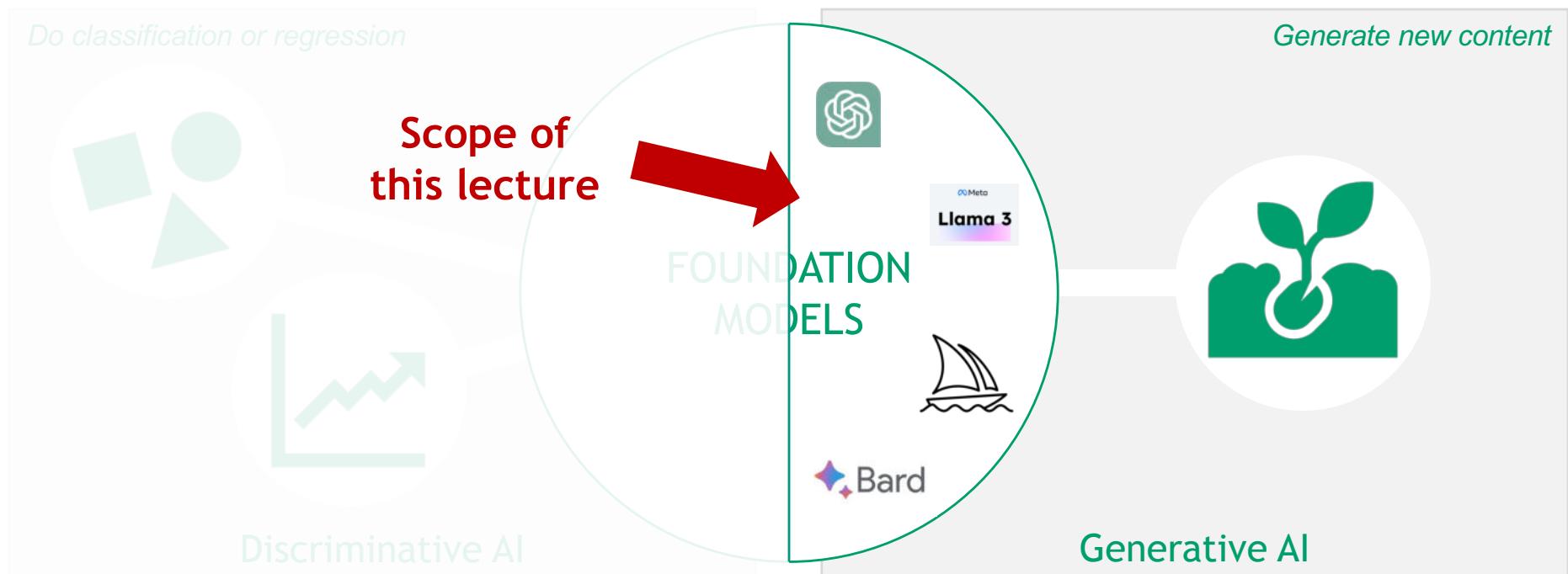
How do Foundation Models relate to Generative AI?

Foundation models are a key component of Generative AI.



How do Foundation Models relate to Generative AI?

Foundation models go well beyond language.





- 1 Introduction
- 2 Multimodality in Data
- 3 Generative Foundation Models
- 4 Research
- 5 Industry

How to deal with multimodality in data?

Images can be described and generated using text prompts.



*"Create a realistic, high-resolution image of **two people standing in a museum**, observing a famous painting [...]. The painting they are observing is inspired by **Vincent van Gogh's 'Self-Portrait with Bandaged Ear'**, characterized by vibrant colors and expressive brush strokes. The two observers are diverse; one is a middle-aged Black woman with shoulder-length hair, wearing a smart-casual outfit, and the other is a young South Asian man dressed in a trendy jacket and jeans [...]"*

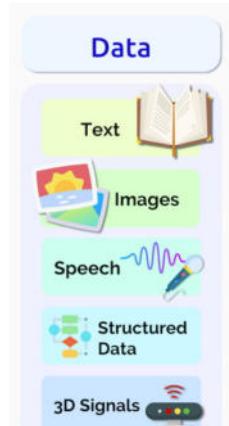


Text and Image are created with ChatGPT and Dall-E of OpenAI

How to deal with multimodality in data?

Textual data is embedded in vectors with semantic relationships.

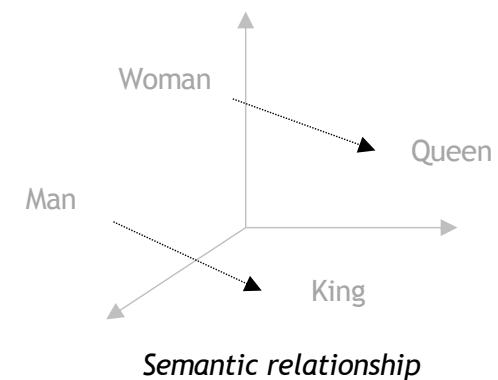
Example



*“two people
standing in a
museum”*

[5]

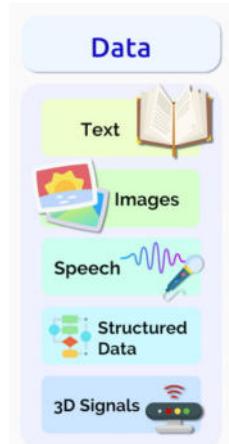
Embedding



How to deal with multimodality in data?

Textual data is embedded in vectors with semantic relationships.

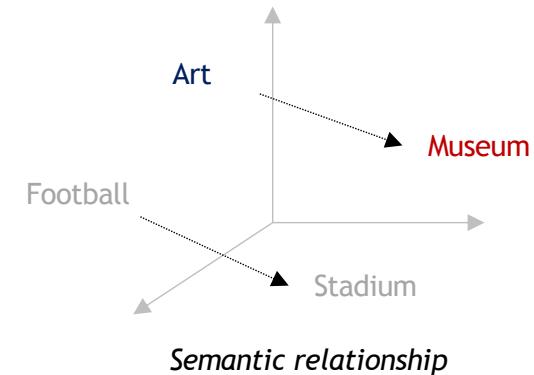
Example



*“two people
standing in a
museum”*

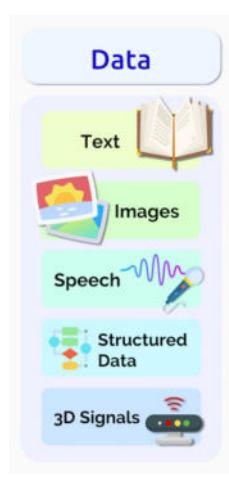
[5]

Embedding



How to deal with multimodality in data?

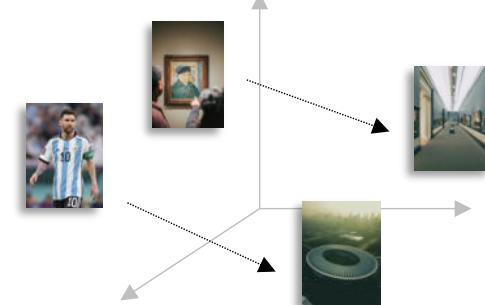
Visual data is embedded in vectors with semantic relationship.



Example



Embedding



How to deal with multimodality in data?
Visual data is embedded in vectors with semantic relationship.

Example



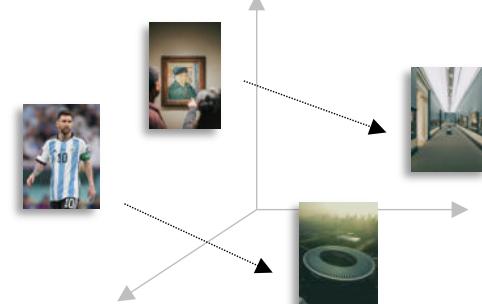
(103,88,67)	(138,118,94)
(138,118,94)	(103,88,67)

(103,88,67)

Image

Pixel RGB Values

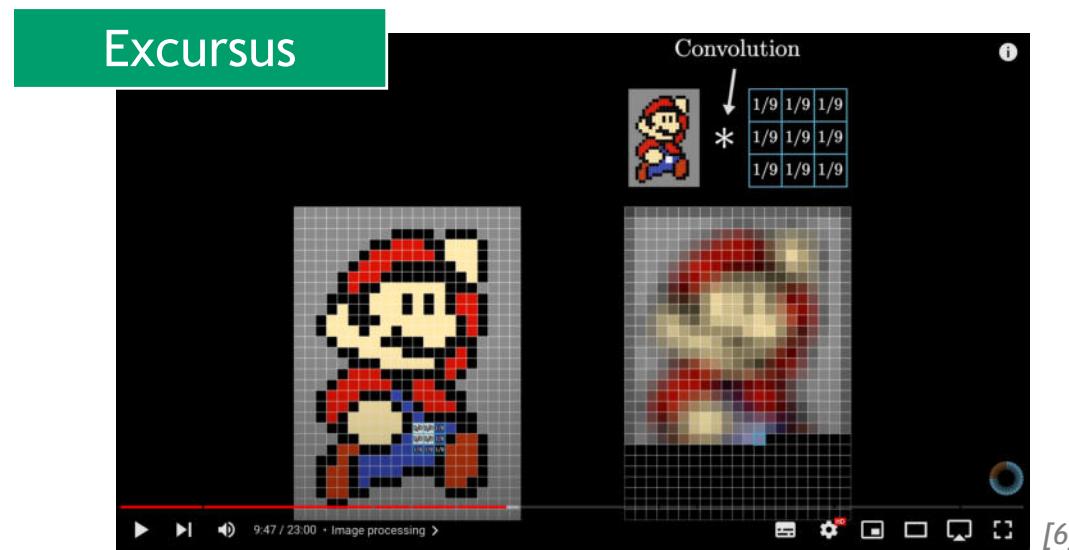
Embedding



Semantic relationship

How to deal with multimodality in data?

Convolution - the basis for computer vision in machine learning.

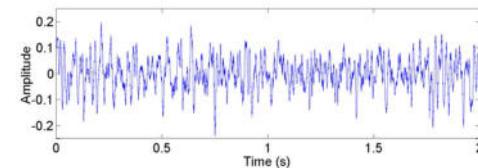
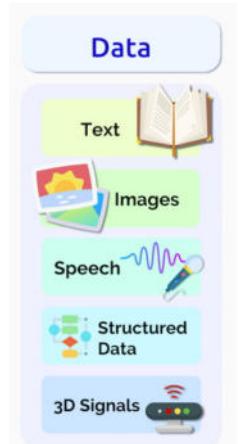


3Blue1Brown, <https://www.youtube.com/watch?v=KuXiwB4LzSA> [6]

How to deal with multimodality in data?

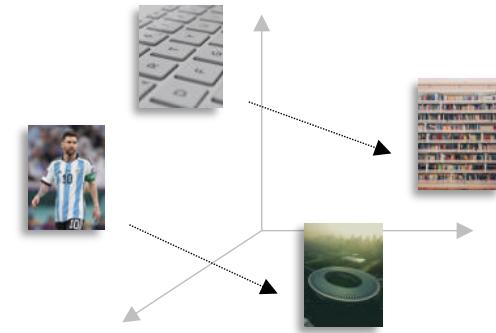
Audio data is embedded in vectors with semantic relationship.

Example



Audio

Embedding

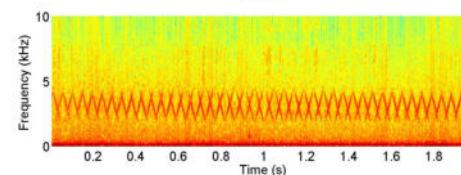
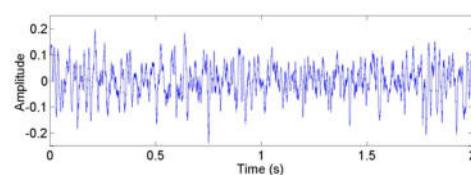


Semantic relationship

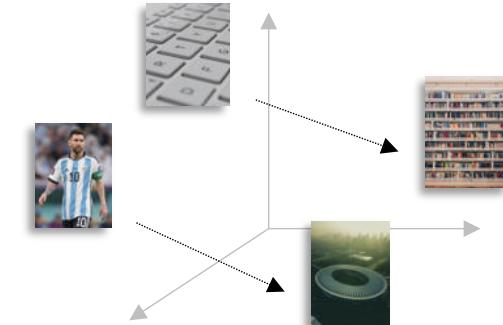
How to deal with multimodality in data?

Audio data is embedded in vectors with semantic relationship.

Example



Spectrogram

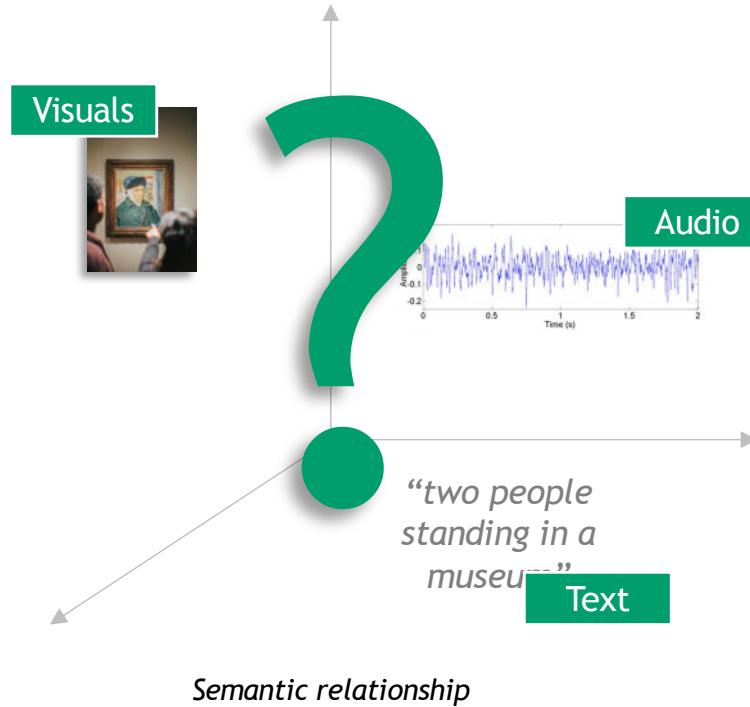


Embedding

Audio

How to deal with multimodality in data?

Embeddings enable semantic relationship of multimodal data.



What are the challenges in multimodal data?
A foundation model can only be as good as the data itself.

“The rising of the sun over the world”



Task

Giovanni Battista Tiepolo
(1696-1770)

Ceiling fresco

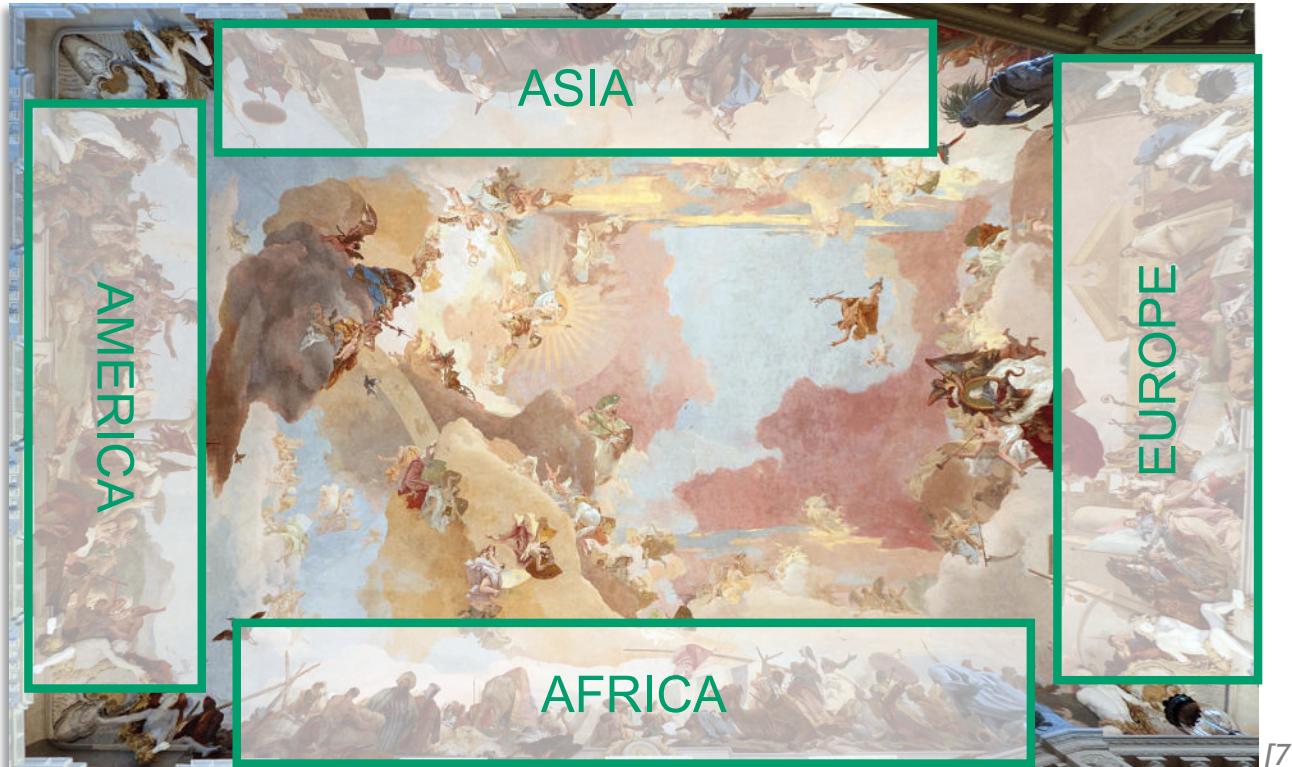
What are the challenges in multimodal data?
A foundation model can only be as good as the data itself.



<https://schloesserblog.bavern.de/heute-vor/giovanni-battista-tiepolo-und-sein-meisterwerk-in-der-residenz-wuerzburg> [7]

What are the challenges in multimodal data?

A foundation model can only be as good as the data itself.



Completeness

[7]

<https://schloesserblog.bavern.de/heute-vor/giovanni-battista-tiepolo-und-sein-meisterwerk-in-der-residenz-wuerzburg> [7]

What are the challenges in multimodal data?

A foundation model can only be as good as the data itself.



Completeness

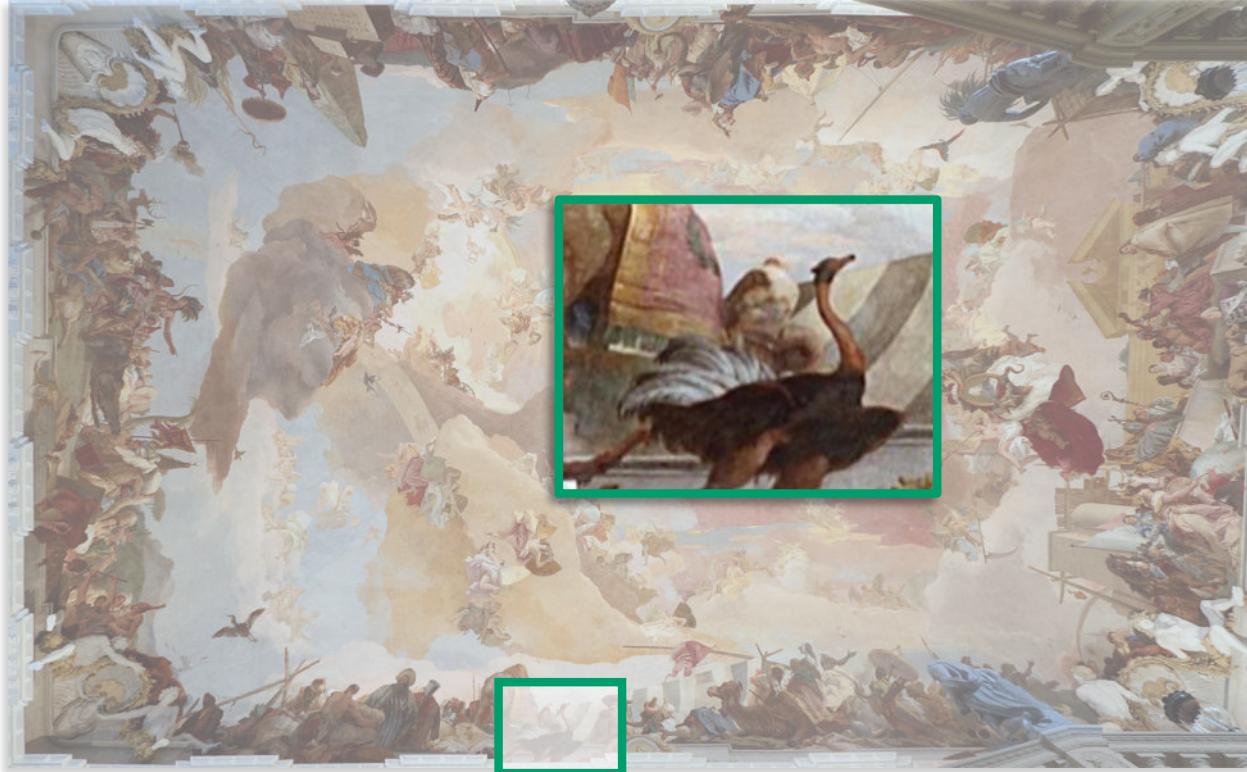
Context

[7]

<https://schloesserblog.bavern.de/heute-vor/giovanni-battista-tiepolo-und-sein-meisterwerk-in-der-residenz-wuerzburg> [7]

What are the challenges in multimodal data?

A foundation model can only be as good as the data itself.



Completeness

Context

Relationship

What are the challenges in multimodal data?

A foundation model can only be as good as the data itself.



Completeness

Context

Relationship

Bias

[7]

<https://schloesserblog.bavern.de/heute-vor/giovanni-battista-tiepolo-und-sein-meisterwerk-in-der-residenz-wuerzburg> [7]

What are the challenges in multimodal data?

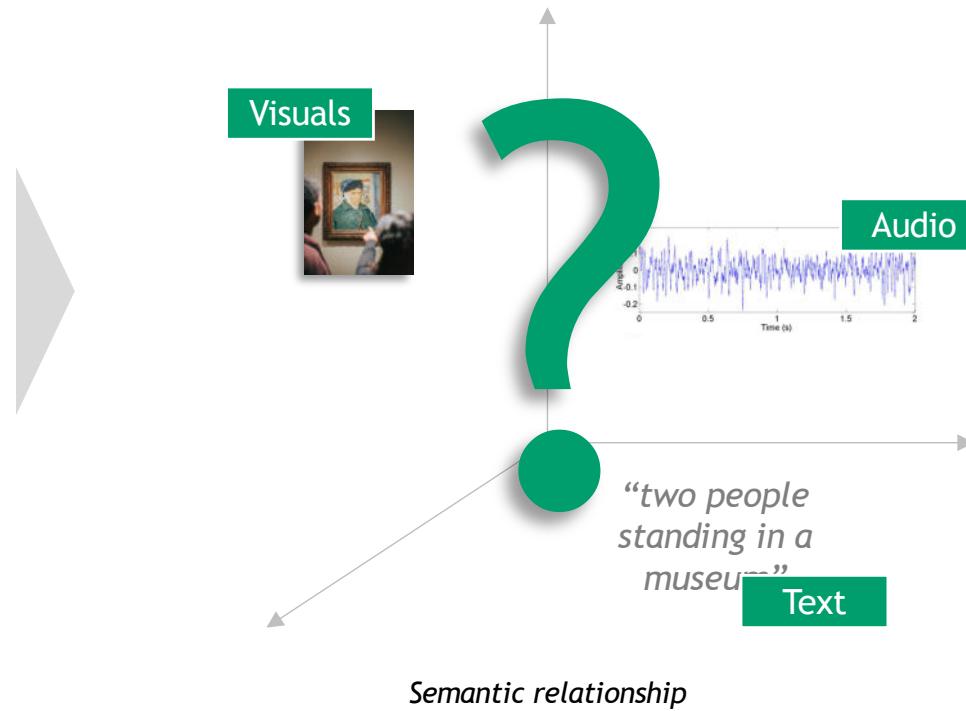
A foundation model can only be as good as the data itself.

Completeness

Context

Relationship

Bias



How to deal with challenges in multimodal data?

Foundation Model outputs a text prompt from most likely next token prediction.



"Create a realistic, high-resolution image of two people standing in a museum, observing a famous painting [...]. The painting they are observing is inspired by Vincent van Gogh's 'Self-Portrait with Bandaged Ear', characterized by vibrant colors and expressive brush strokes. The two observers are diverse; one is a middle-aged Black woman with shoulder-length hair, wearing a smart-casual outfit, and the other is a young South Asian man dressed in a trendy jacket and jeans [...]"



Text and Image are created with ChatGPT and Dall-E of OpenAI

How to deal with challenges in multimodal data?

Foundation Model mirrors “information” seen in training data.



"Create a realistic, high-resolution image of two people standing in a museum, observing a famous painting [...]. The painting they are observing is inspired by Vincent van Gogh's 'Self-Portrait with Bandaged Ear', characterized by vibrant colors and expressive brush strokes. The two observers are diverse; one is a middle-aged Black woman with shoulder-length hair, wearing a smart-casual outfit, and the other is a young South Asian man dressed in a trendy jacket and jeans [...]"



Text and Image are created with ChatGPT and Dall-E of OpenAI

How to deal with challenges in multimodal data?

Foundation Model makes assumptions based on frequency of occurrences.



*"Create a realistic, high-resolution image of two people standing **in a museum**, observing a famous painting [...]. The painting they are observing is **inspired by** Vincent van Gogh's 'Self-Portrait with Bandaged Ear', characterized by vibrant colors and expressive brush strokes. The **two observers are diverse; one is a middle-aged Black woman with shoulder-length hair, wearing a smart-casual outfit, and the other is a young South Asian man dressed in a trendy jacket and jeans [...]"***



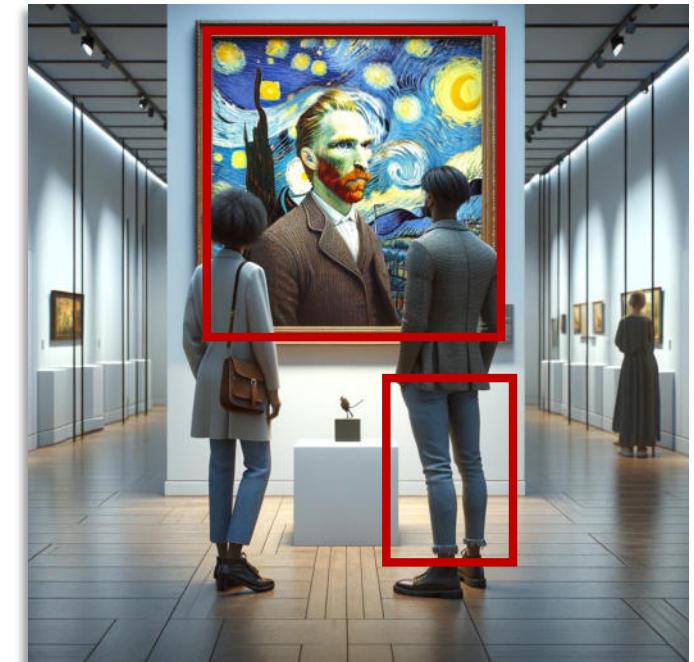
Text and Image are created with ChatGPT and Dall-E of OpenAI

How to deal with challenges in multimodal data?

Foundation Model “mutates and recombines” the data - partly adding new content.

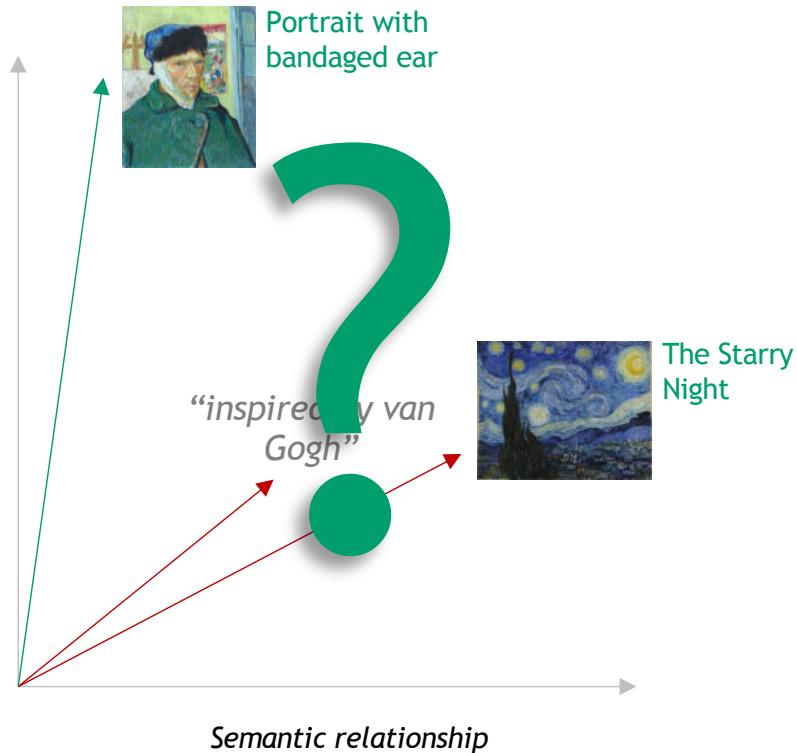


"Create a realistic, high-resolution image of two people standing in a museum, observing a famous painting [...]. The painting they are observing is inspired by Vincent van Gogh's 'Self-Portrait with Bandaged Ear', characterized by vibrant colors and expressive brush strokes. The two observers are diverse; one is a middle-aged Black woman with shoulder-length hair, wearing a smart-casual outfit, and the other is a young South Asian man dressed in a trendy jacket and jeans [...]"



Text and Image are created with ChatGPT and Dall-E of OpenAI

How to deal with challenges in multimodal data? Data (embedding) is the driver of foundation models.

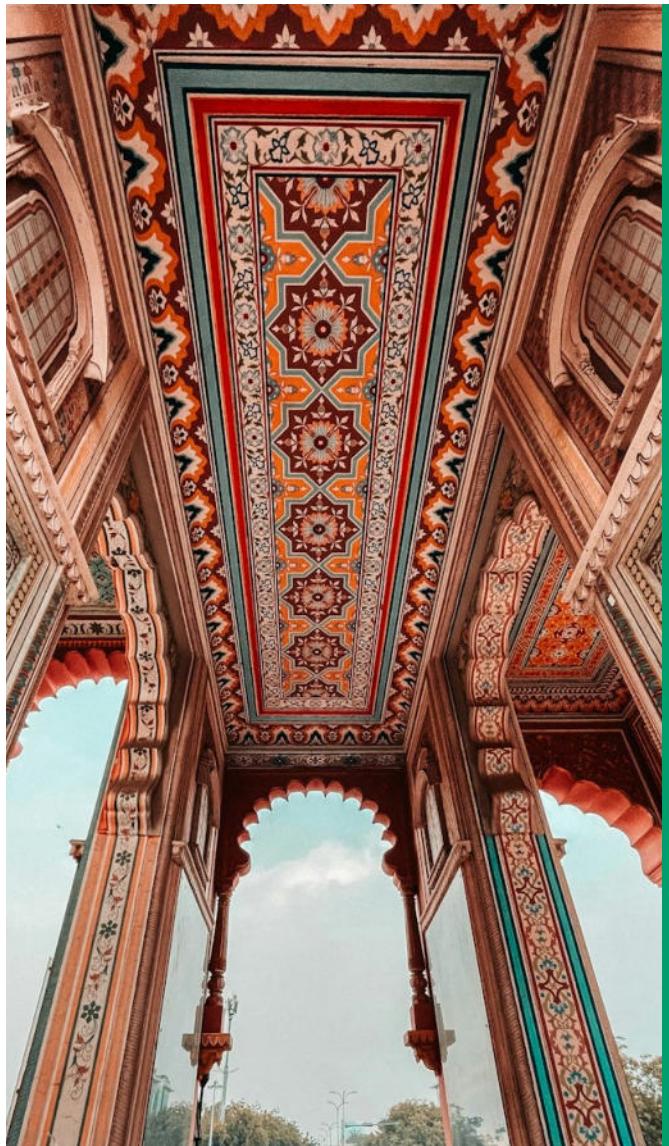


Text and Image are created with ChatGPT and Dall-E of OpenAI

A perspective on Foundation Models

What power arises from data creation and curation?

Discussion



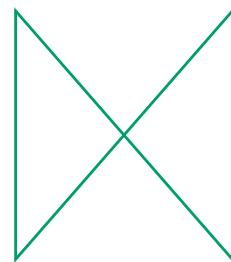
- 1 Introduction
- 2 Multimodality in Data
- 3 Generative Foundation Models
- 4 Research
- 5 Industry

Generative Foundation Models

This lecture focuses on visual and text data.



Foundation
Model



Allegory of the Cave

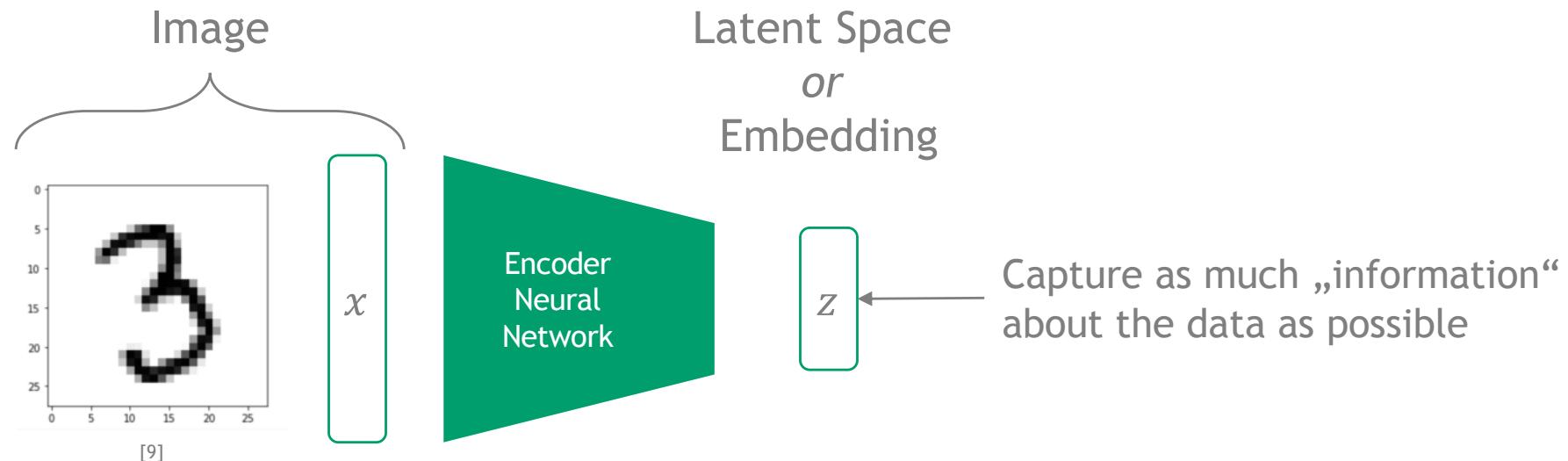
Can we learn the true explanatory factors from solely observed data?



An Illustration of The Allegory of the Cave, from Plato's Republic.jpg [8]

Autoencoder

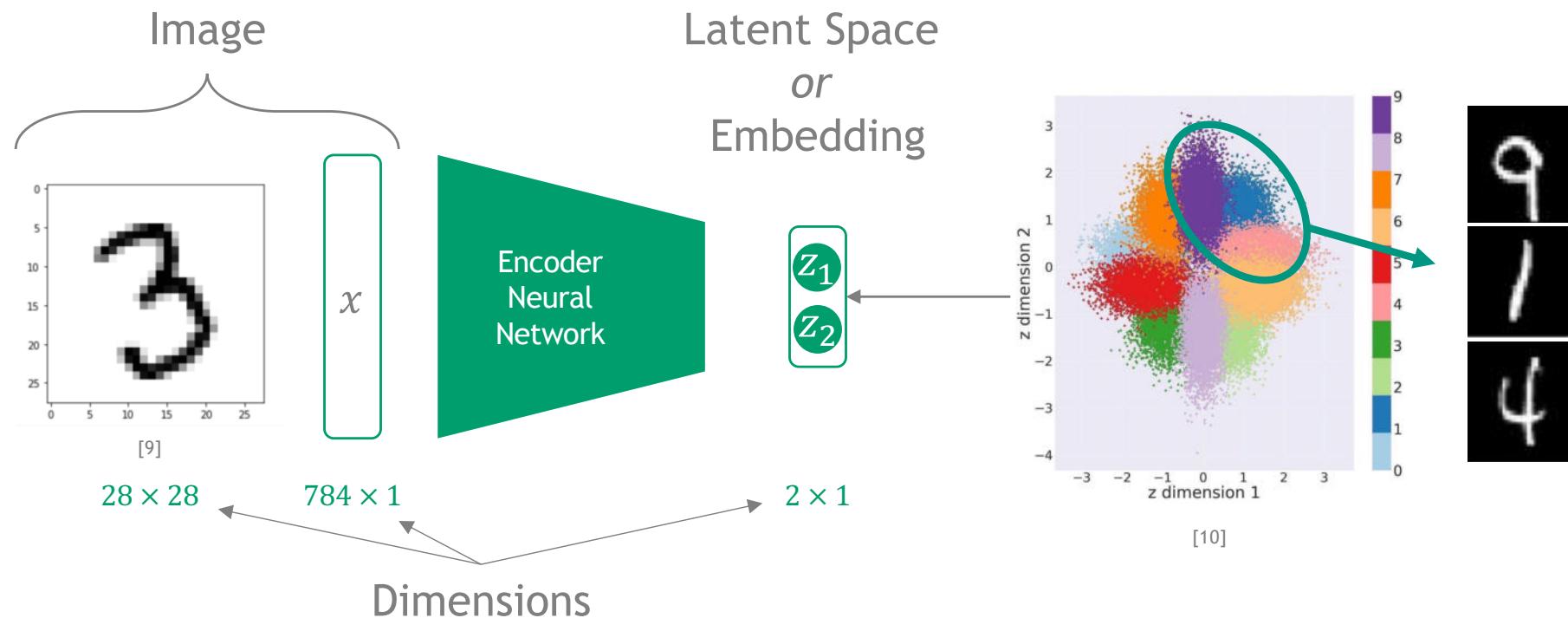
Autoencoding is a form of compression.



30,000 sample in the MNIST dataset. [9]

Autoencoder

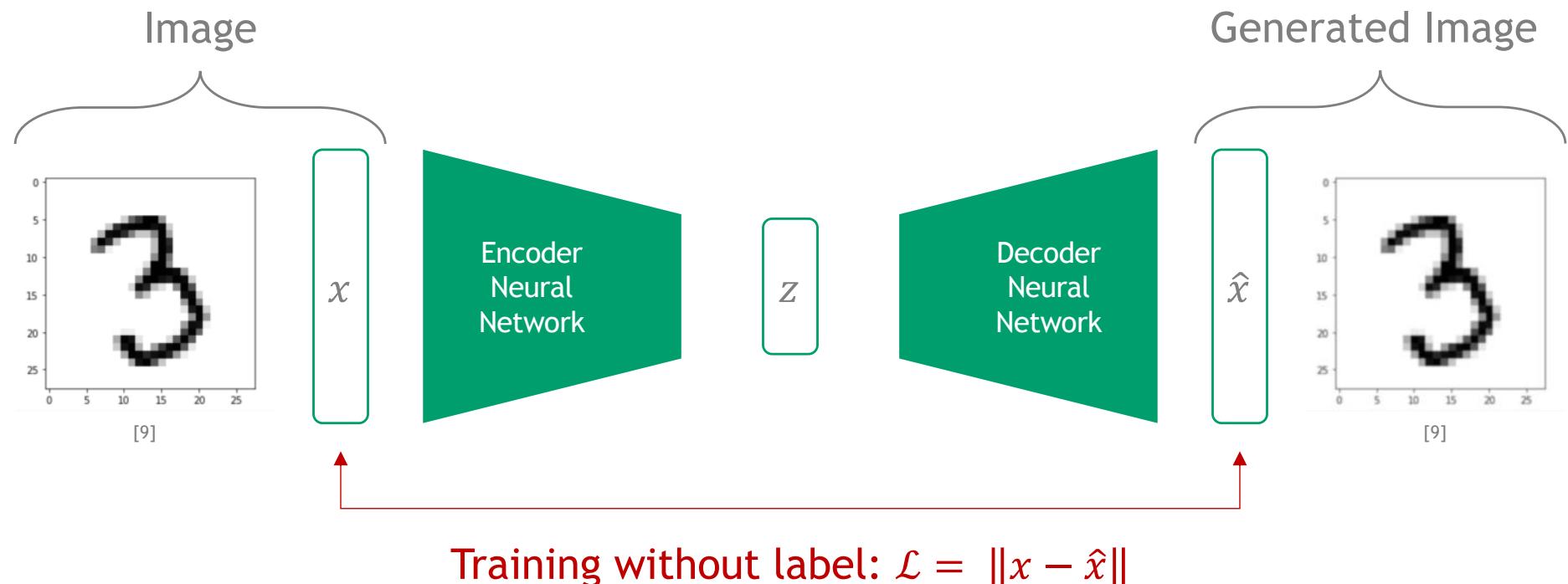
The encoder learns mapping from the data, x , to a low-dimensional space, z .



30,000 sample in the MNIST dataset. [9]
 Mundt et al., (2022). Unified Probabilistic Deep Continual Learning through Generative Replay and Open Set Recognition [10]

Autoencoder

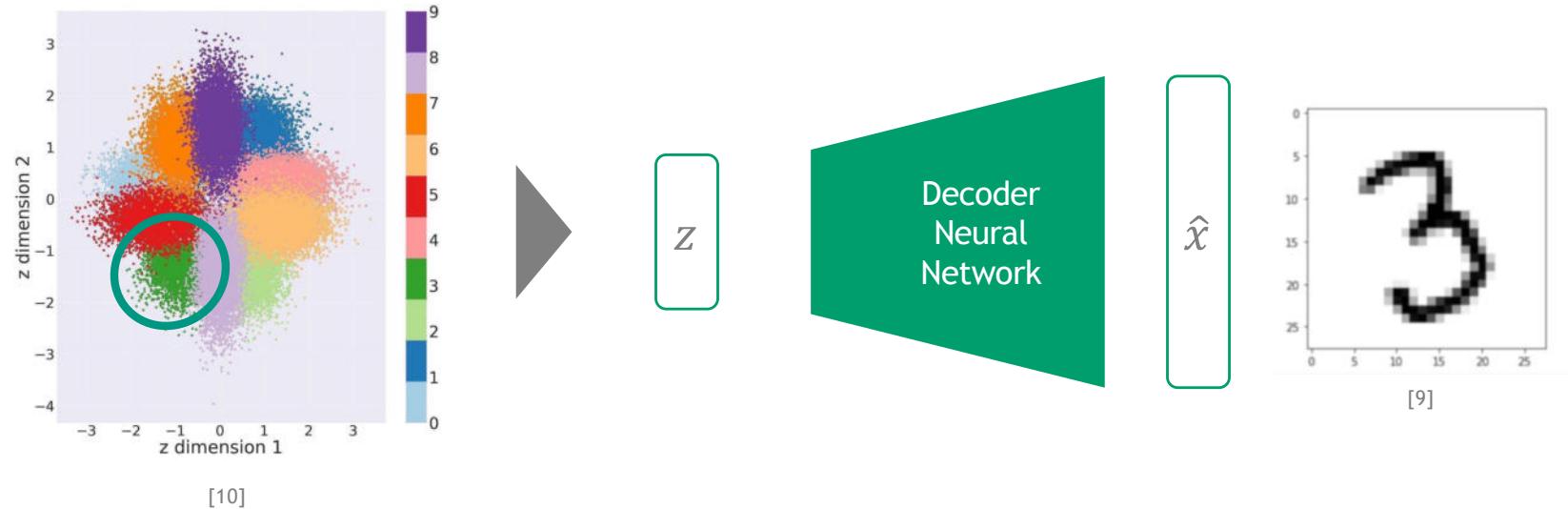
Train the model by reconstructing the original data.



30,000 sample in the MNIST dataset. [9]

Autoencoder

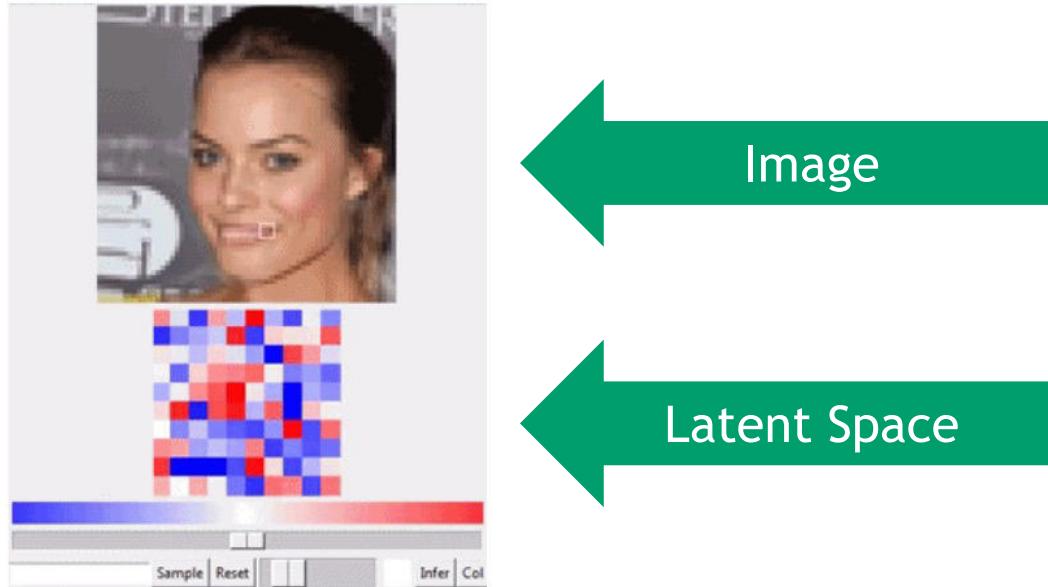
Generate new samples from the latent space.



30,000 sample in the MNIST dataset. [9]
Mundt et al., (2022). Unified Probabilistic Deep Continual Learning through Generative Replay and Open Set Recognition [10]

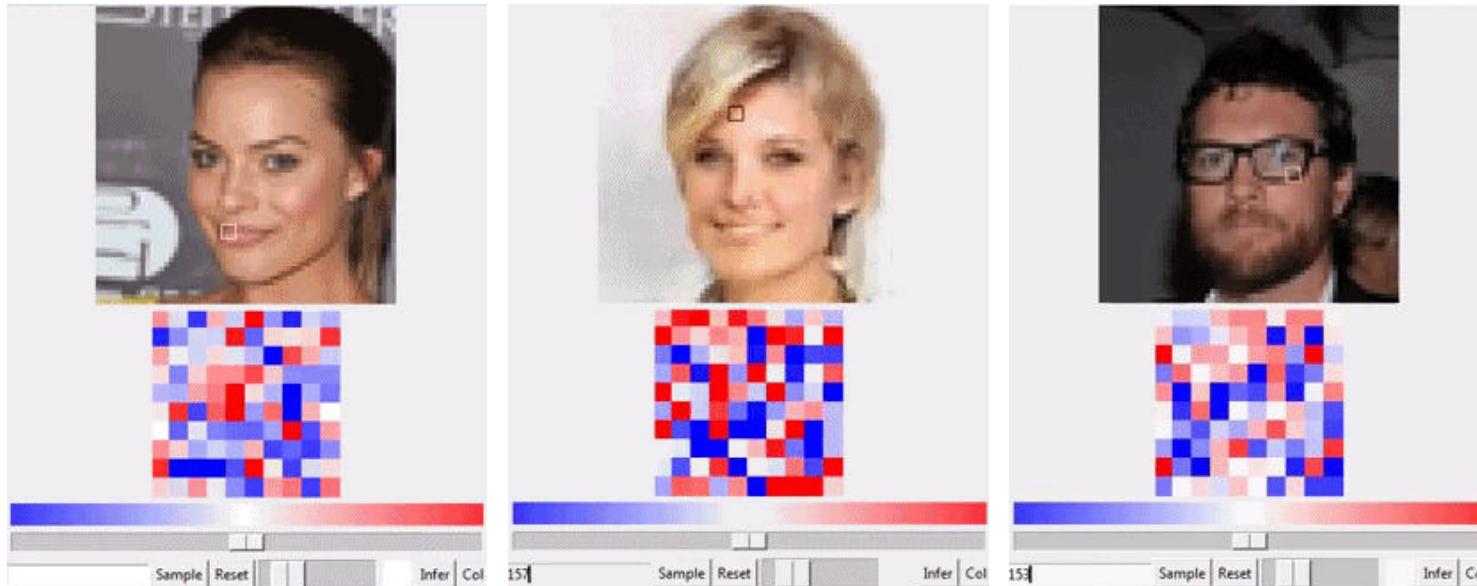
Autoencoder

The latent space can compress representation of world.



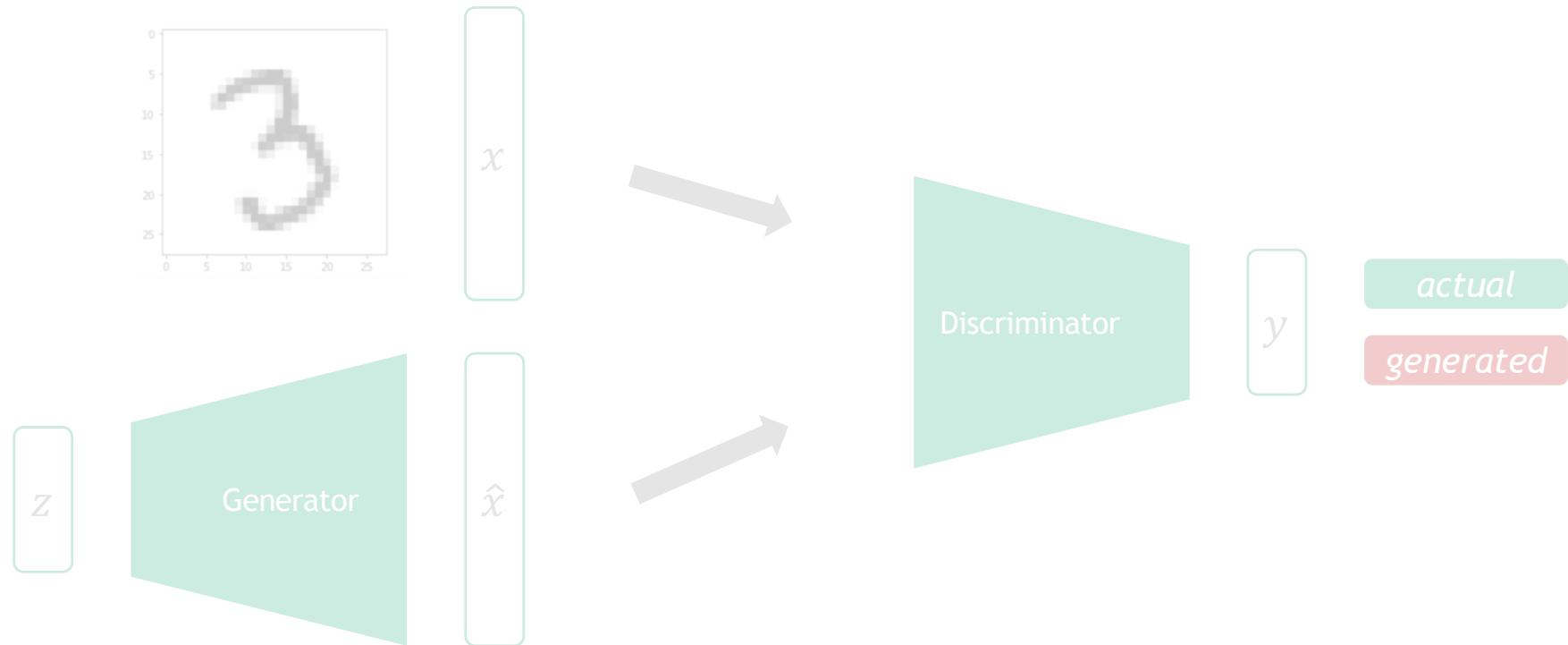
Autoencoder

The latent space can compress representation of world.



Generative Adversarial Network (GAN)

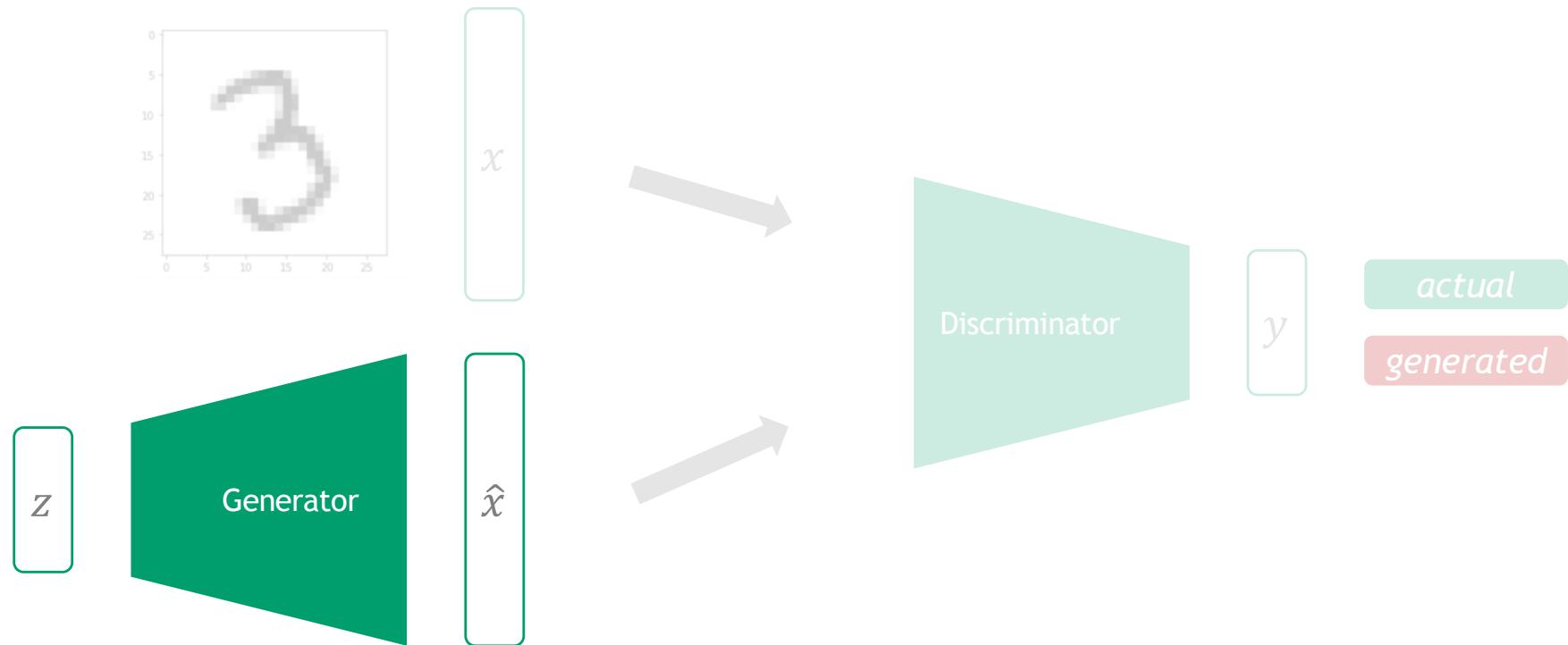
Just sample to generate new instances using competing neural networks.



Goodfellow et al., (2014). Generative Adversarial Nets

Generative Adversarial Network (GAN)

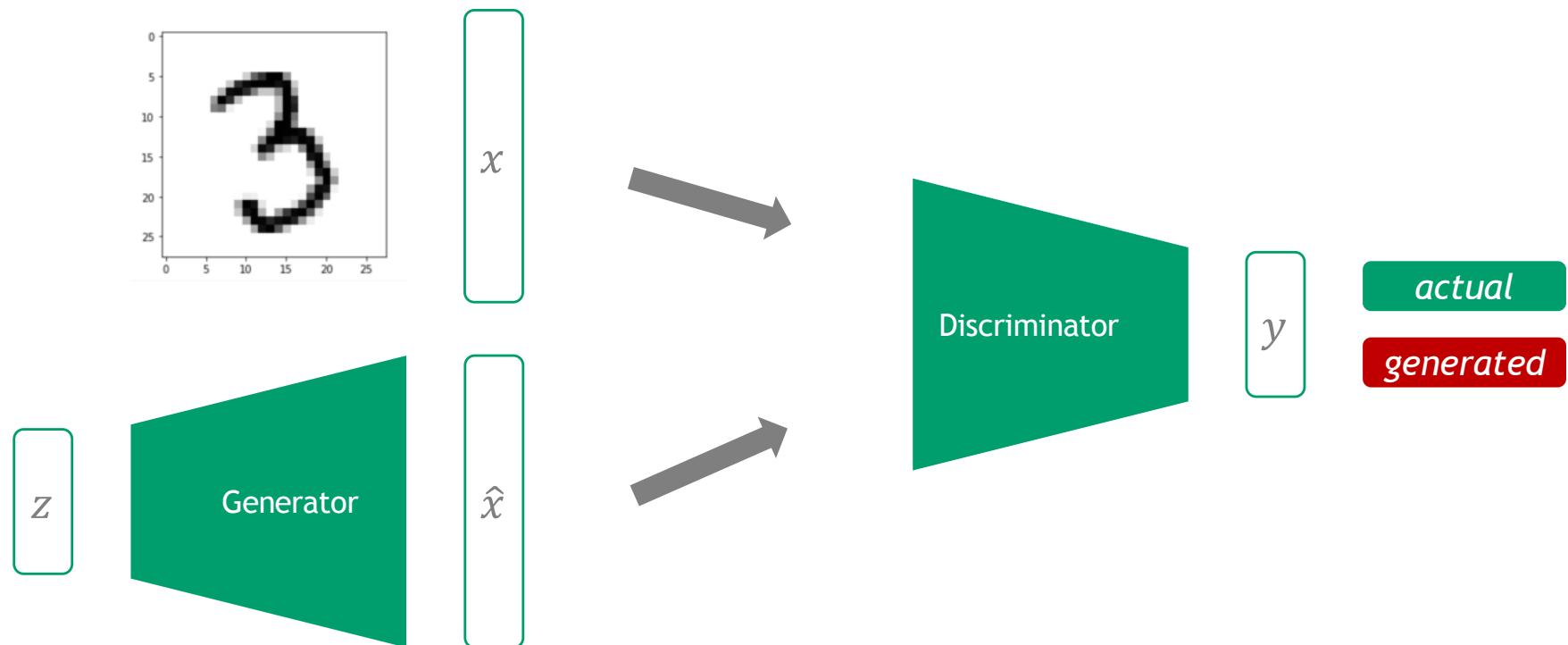
Generator turns noise into an imitation of the data.



Goodfellow et al., (2014). Generative Adversarial Nets

Generative Adversarial Network (GAN)

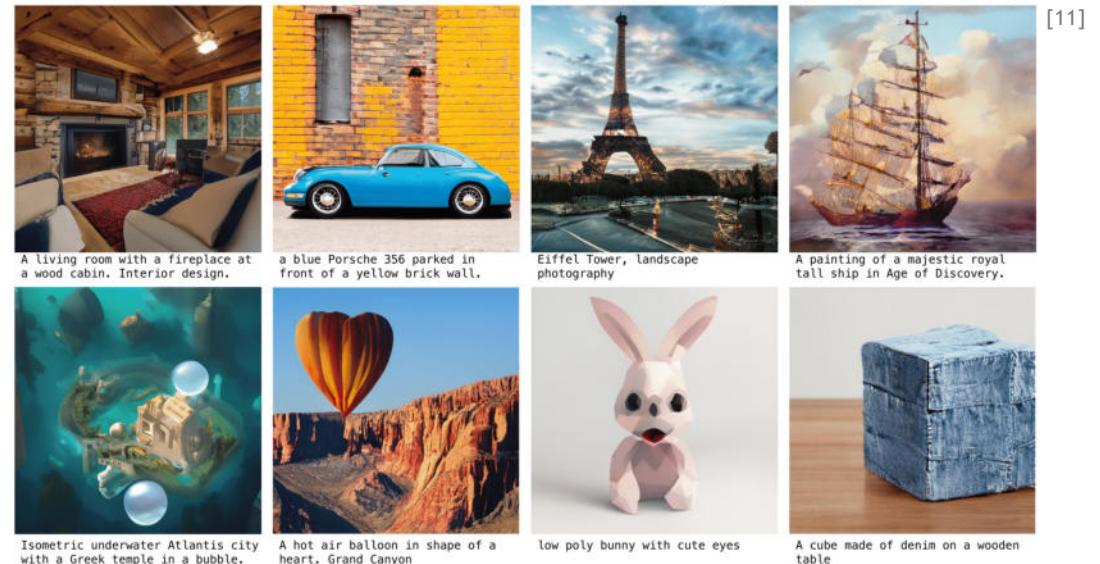
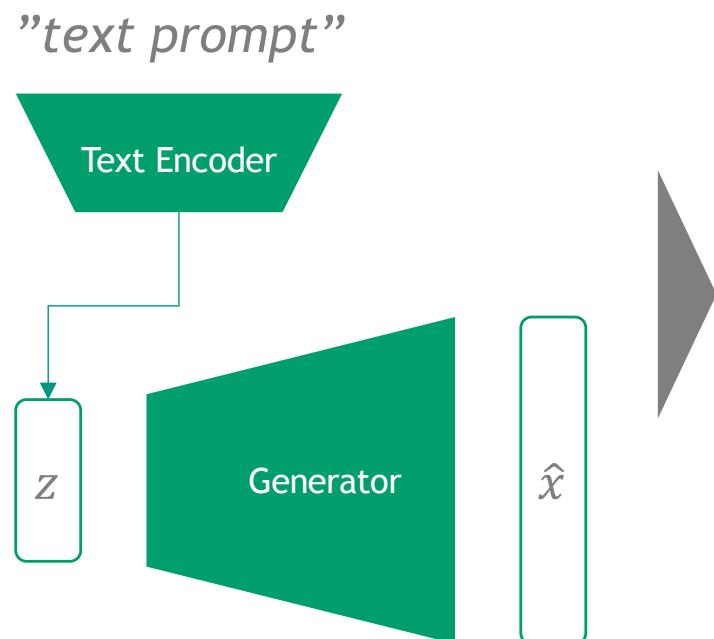
Discriminator tries to identify actual data from fakes created by the generator.



Goodfellow et al., (2014). Generative Adversarial Nets

Generative Adversarial Network (GAN)

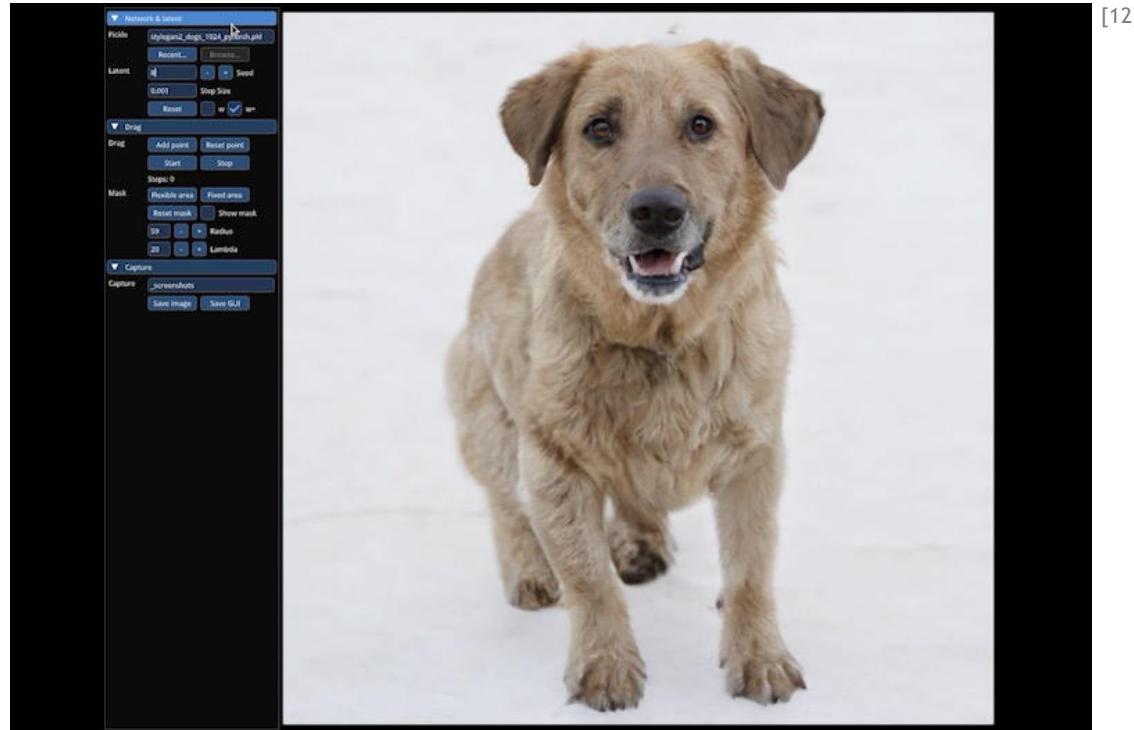
Image generation using encoded text prompts.



Kang et al., (2023). Scaling up GANs for Text-to-Image Synthesis [11]

Generative Adversarial Network (GAN)

Image generation by conditioning the latent space.

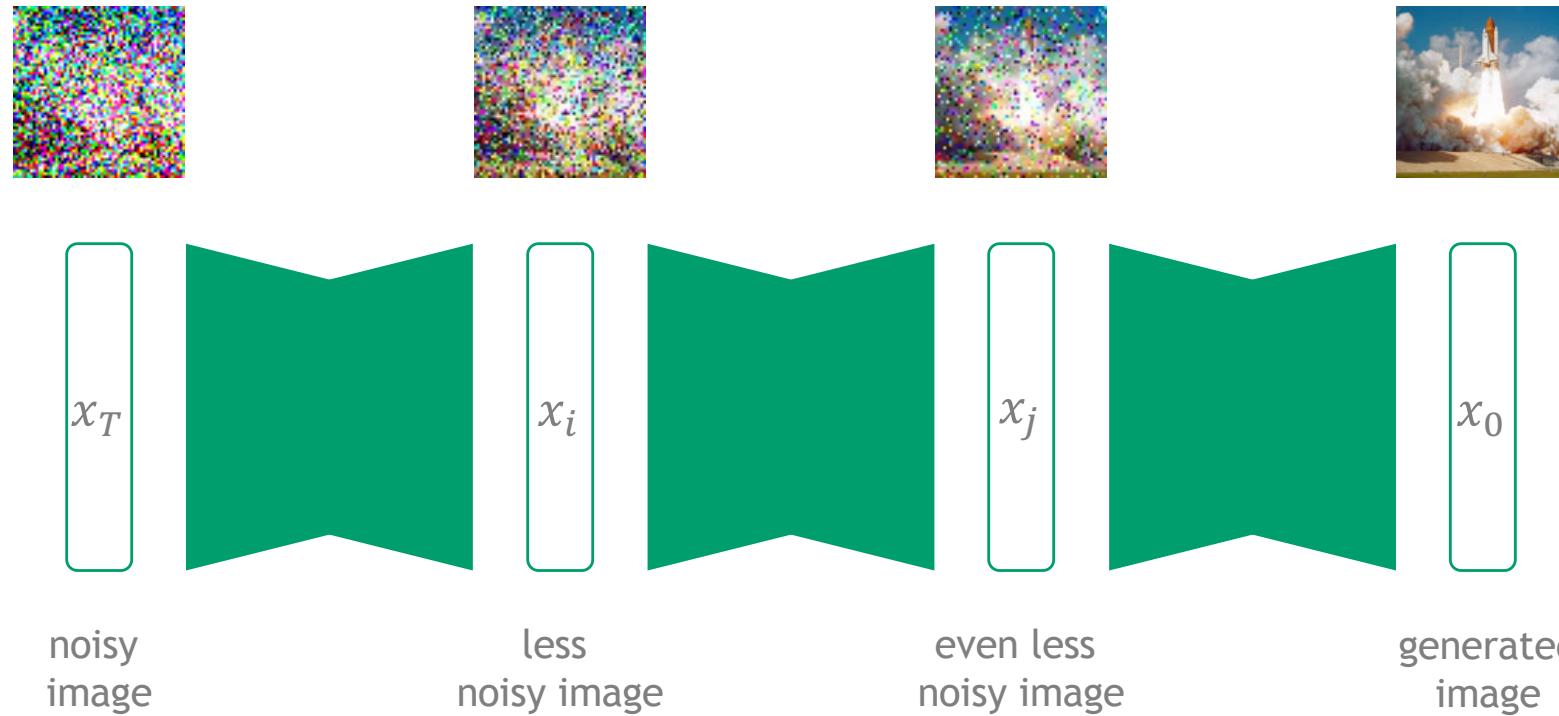


[12]

Pan et al., (2023). Drag Your GAN: Interactive Point-based Manipulation on the Generative Image Manifold [12]

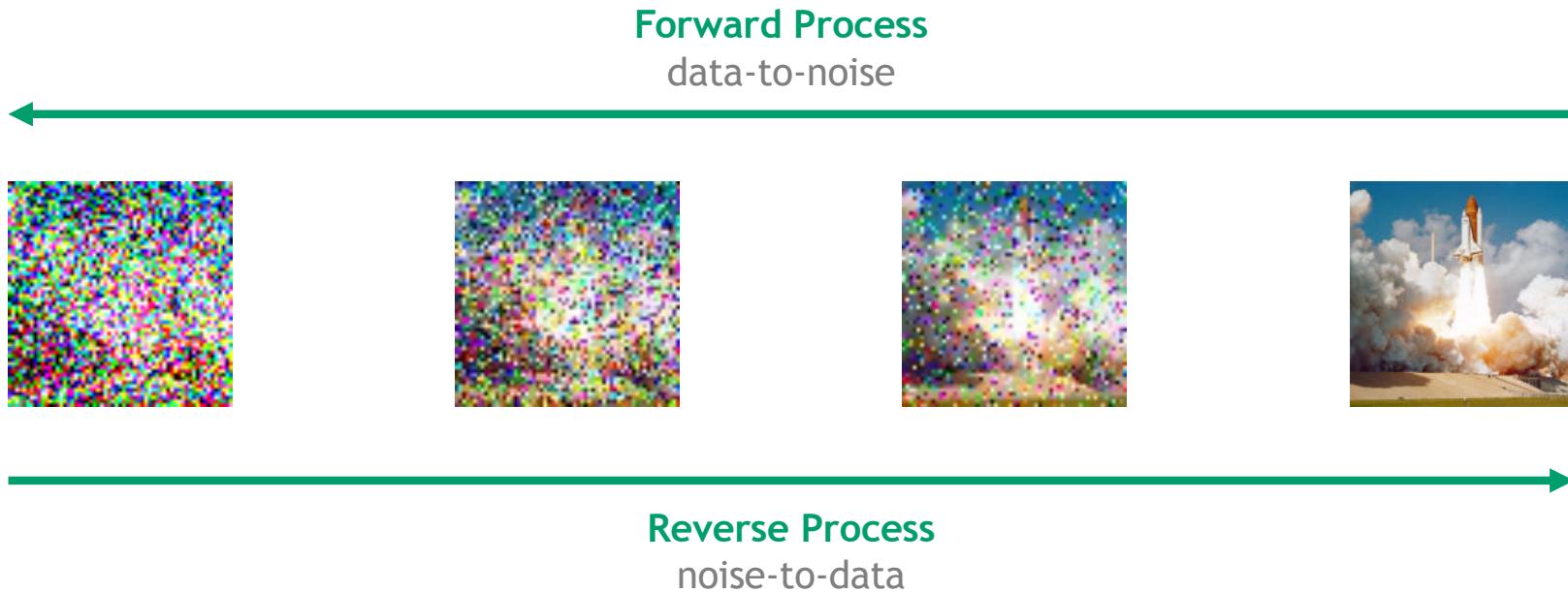
Diffusion Model

Generating images iteratively by repeatedly removing noise.



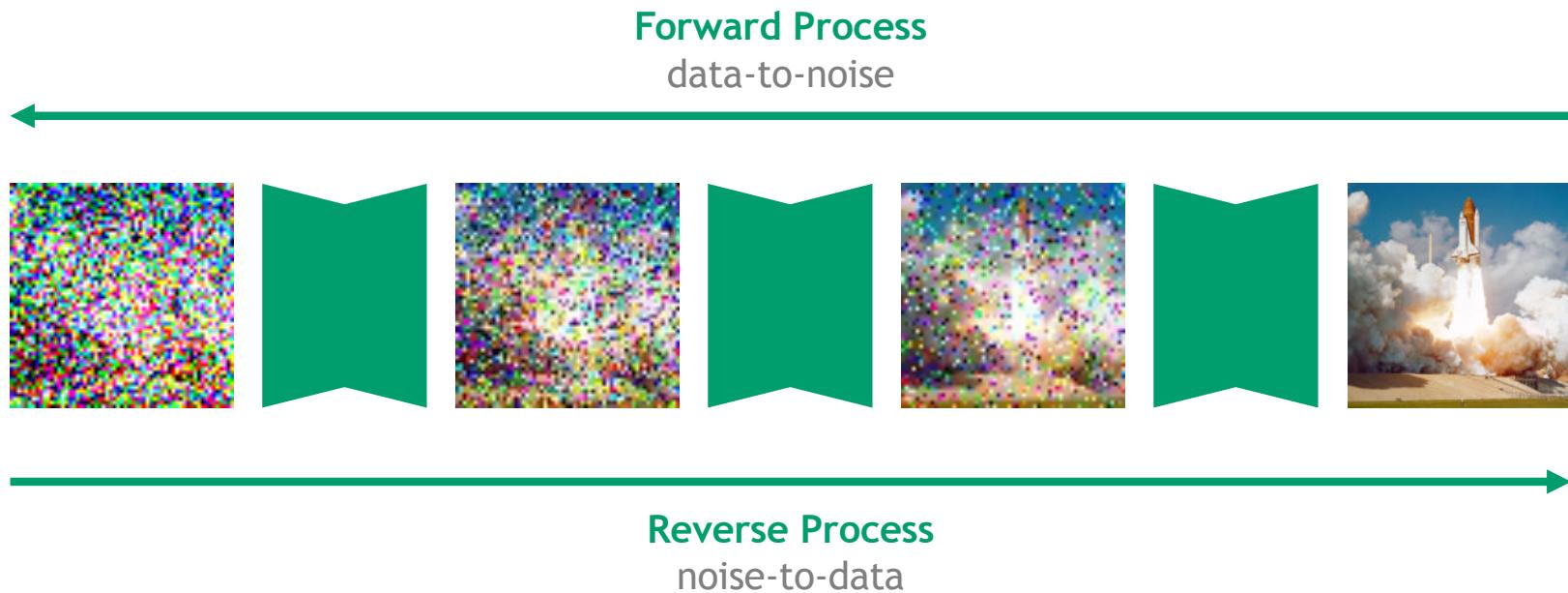
Diffusion Model

The concept is composed of diffusion and denoising.



Diffusion Model

The model predicts the noise in each iteration.



Diffusion Model

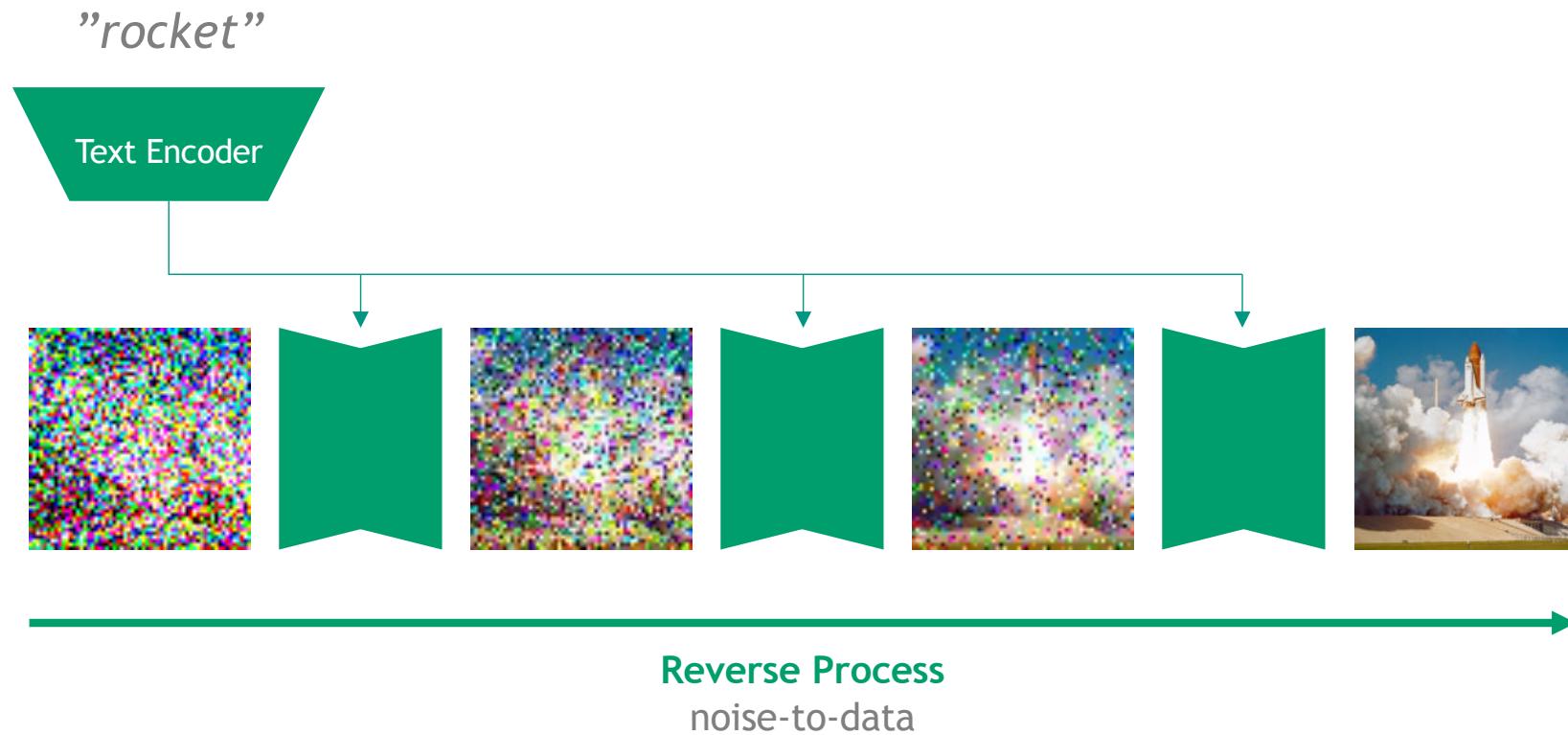
Sampling brand new generations from pure Gaussian noise.



Reverse Process
noise-to-data

Diffusion Model

Generation can be conditioned with (encoded) text prompts.



Diffusion Model

Powerful text-to-video synthesis.



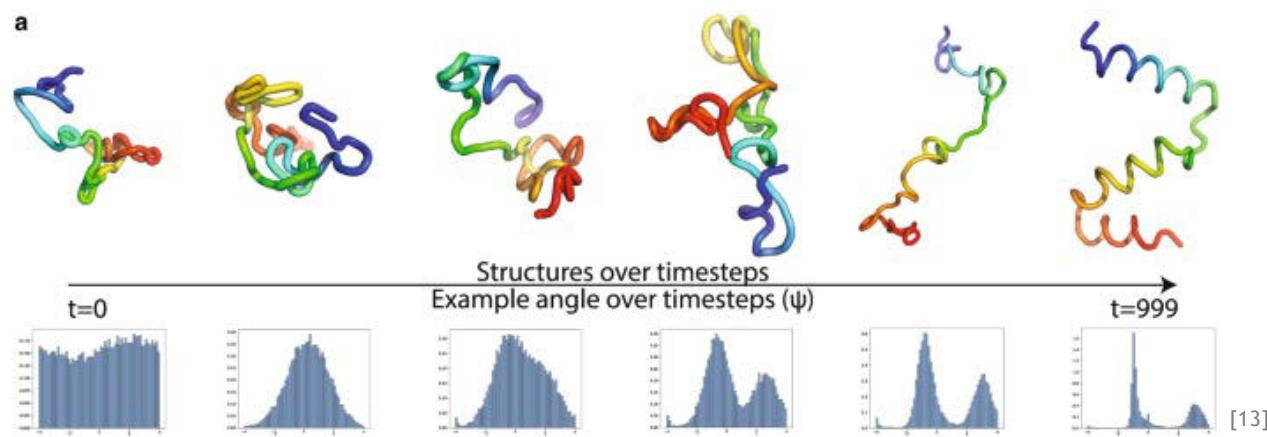
A large orange octopus is seen resting on the bottom of the ocean floor, blending in with the sandy and rocky terrain. Its tentacles are spread out around its body, and its eyes are closed. The octopus is unaware of a king crab that is crawling towards it from behind a rock, its claws raised and ready to attack. The crab is brown and spiny, with long legs and antennae. The scene is captured from a wide angle, showing the vastness and depth of the ocean. The water is clear and blue, with rays of sunlight filtering through. The shot is sharp and crisp, with a high dynamic range. The octopus and the crab are in focus, while the background is slightly blurred, creating a depth of field effect.



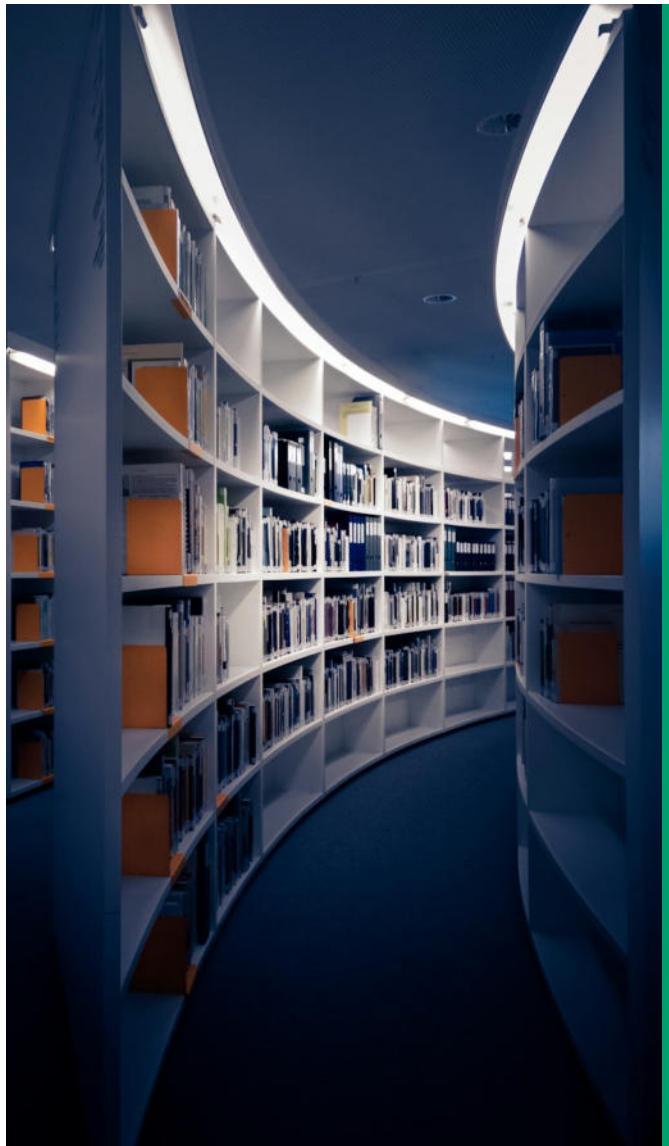
Videos from Sora Model of OpenAI

Diffusion Model

Powerful ~~text-to-image~~ synthesis.



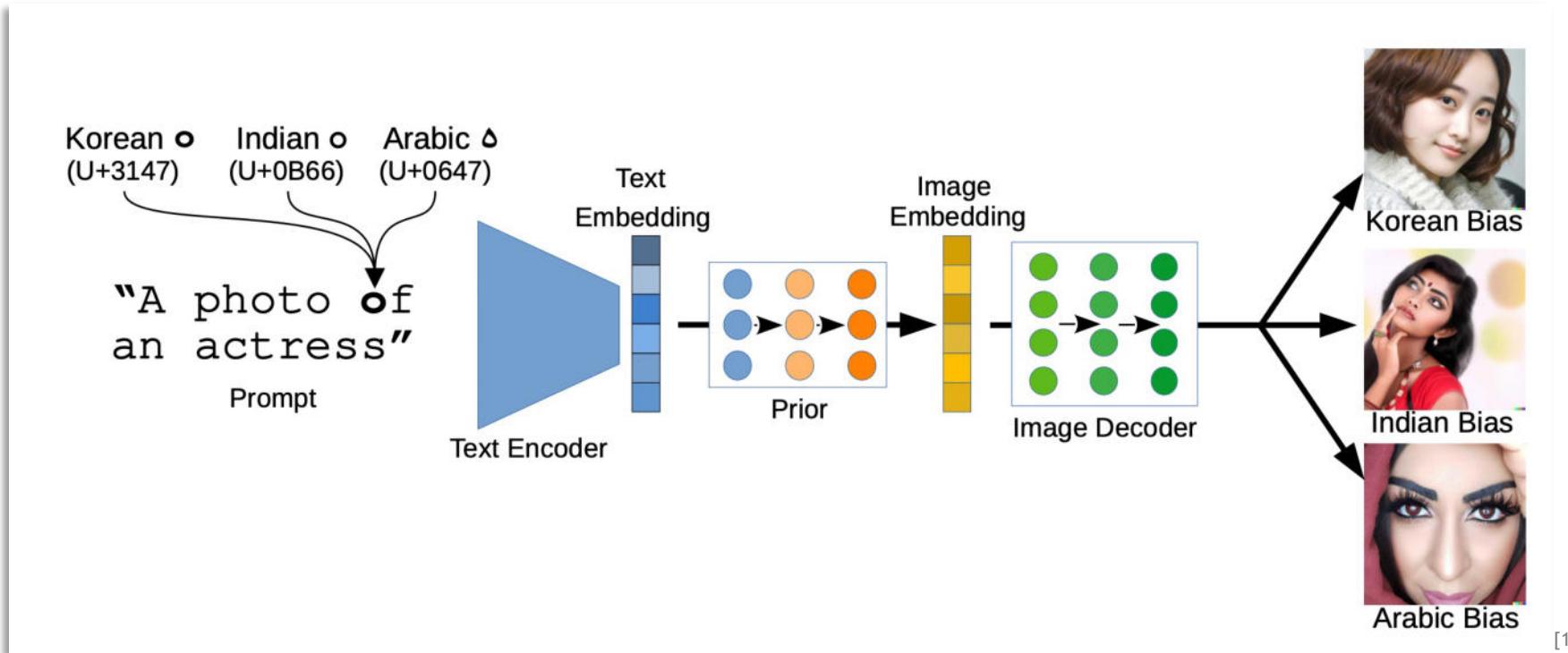
Wu et al., (2023). Protein structure generation via folding diffusion [13]



- 1 Introduction
- 2 Multimodality in Data
- 3 Generative Foundation Models
- 4 Research
- 5 Industry

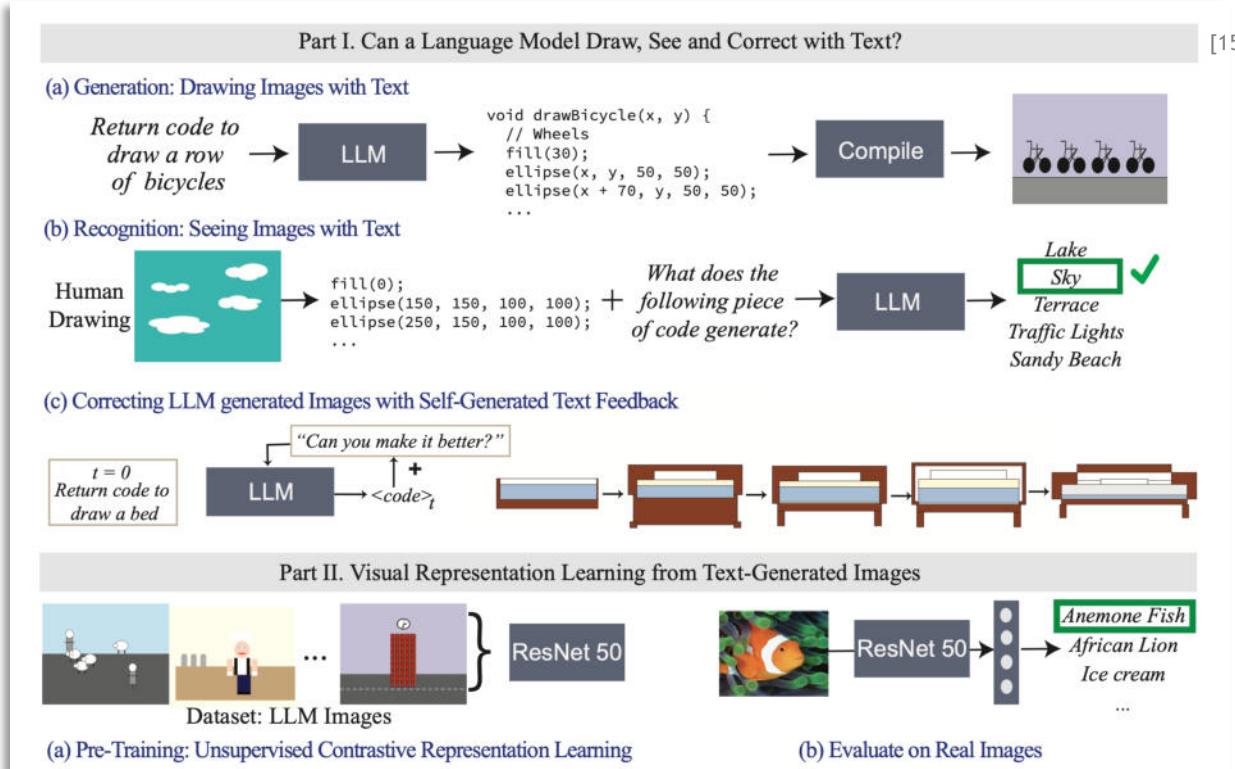
Do we understand the internal function of multimodal models?

Homoglyphs induce stereotypes from cultural circles.



What does it mean to understand the visual concept of a frog?

Language models know a lot about the visual world,



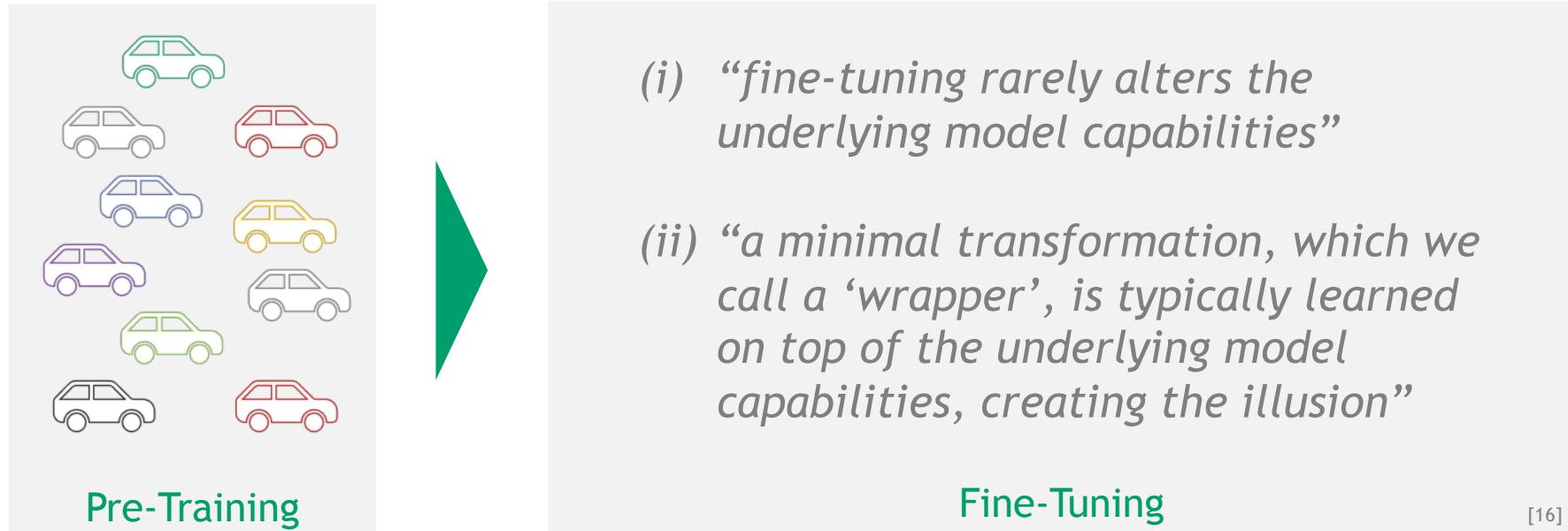
What happens in Fine-Tuning?

Adaptions are trained on top of the pre-training.



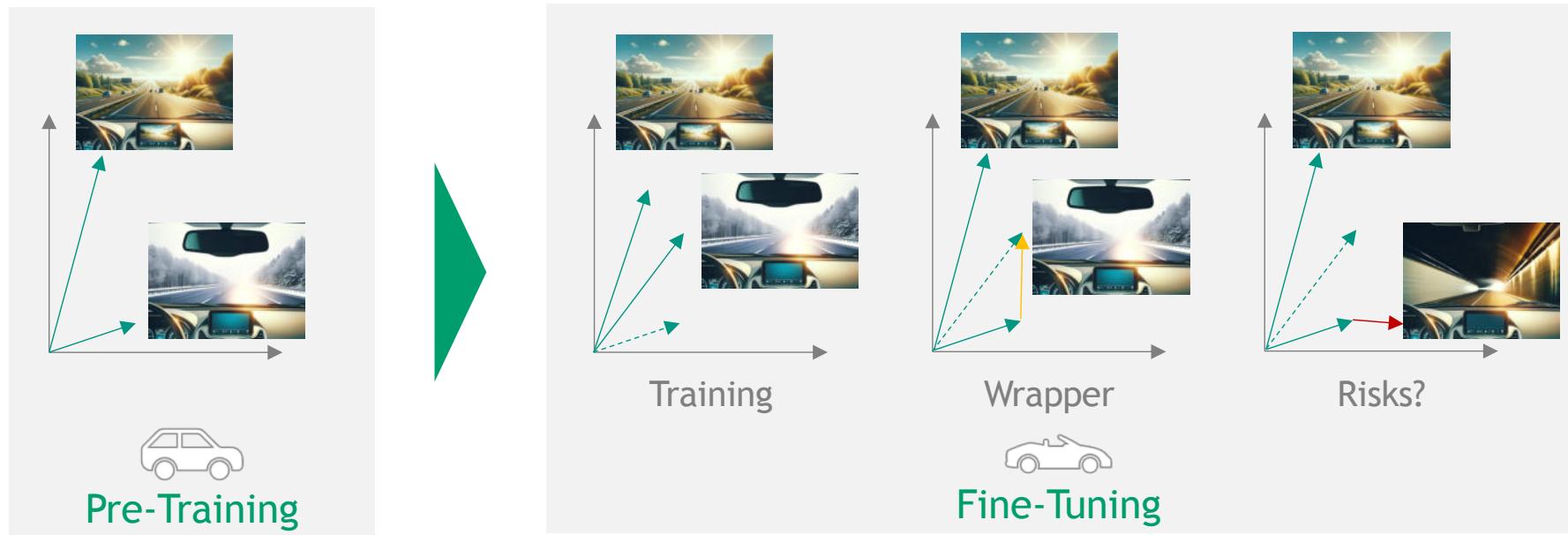
What happens in Fine-Tuning?

The adaption just wraps the underlying representation.



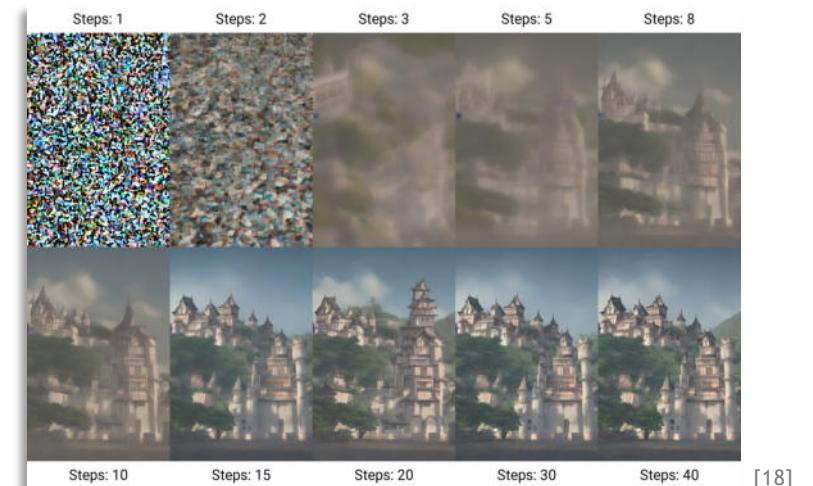
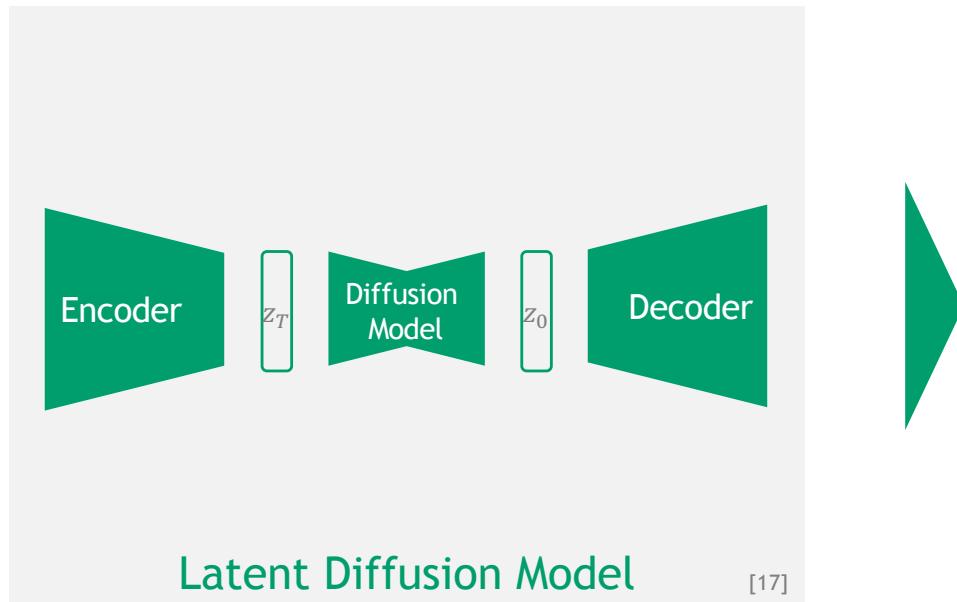
What happens in Fine-Tuning?

Adaptions are trained on top of the pre-training.



Stable Diffusion

Connecting Autoencoders and Diffusion Models.

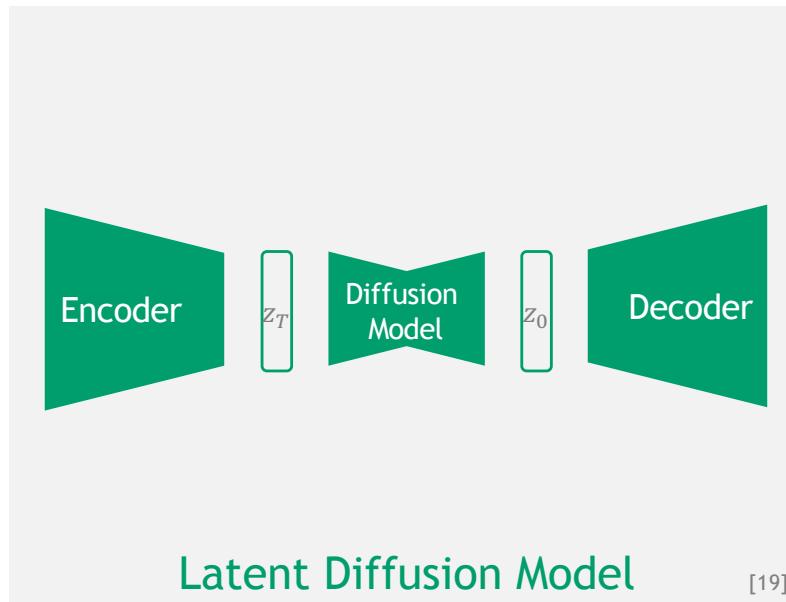


Public Model

Wikipedia, User Bentisquare. [18]
Rombach et al., (2022). High-Resolution Image Synthesis with Latent Diffusion Models. [17]

Stable Diffusion

Approaching Open-Source AI by publishing model & data.



PROJECTS

DATASETS

LAION-400M	Formerly known as crawling@home (C@H), an openly accessible 400M image-text-pair dataset.
LAION5B	A dataset consisting of 5.85 billion CLIP-filtered image-text pairs, featuring several nearest neighbor indices, an improved web-interface for exploration and subset generation, and detection scores for watermark, NSFW, and toxic content detection.
Laion-coco	600M captions generated using BLIP from Laion2B-en.
Laion translated	3B translated samples from Laion5B.

[19]

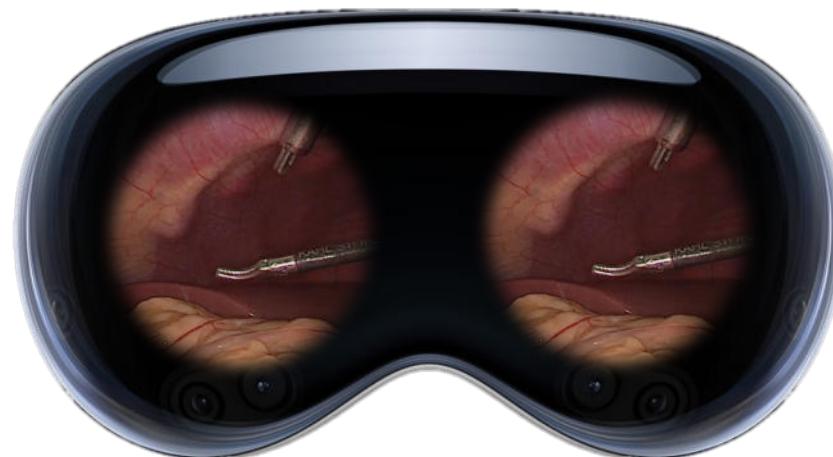
Public Training Data

LAION gemeinnütziger e.V. [19]
Rombach et al., (2022). High-Resolution Image Synthesis with Latent Diffusion Models. [17]

How can Generative AI enhance surgical training?
Use diffusion models to create generated surgical environments.



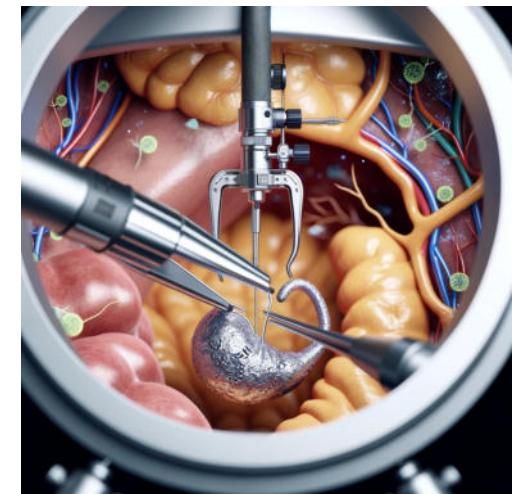
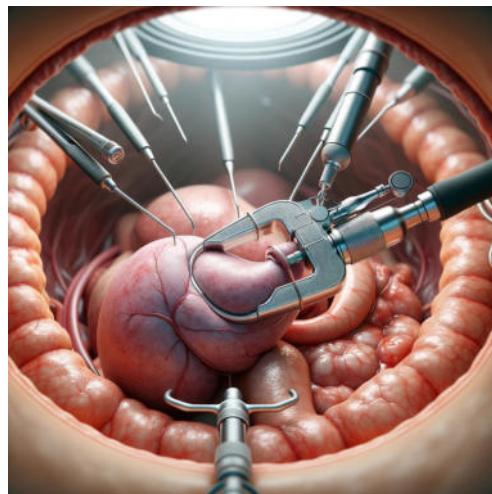
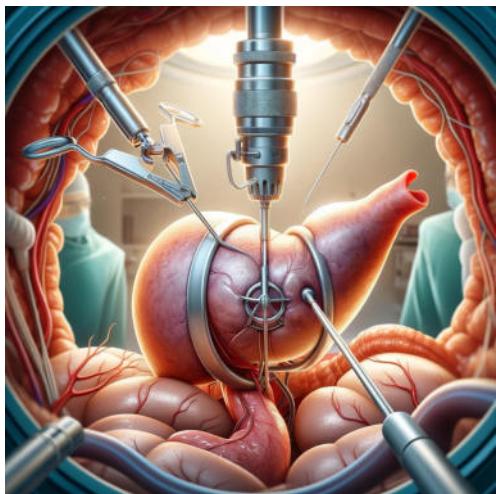
Digital Surgery Lab



Surgical Simulations

<https://www.experimental-surgery.de>

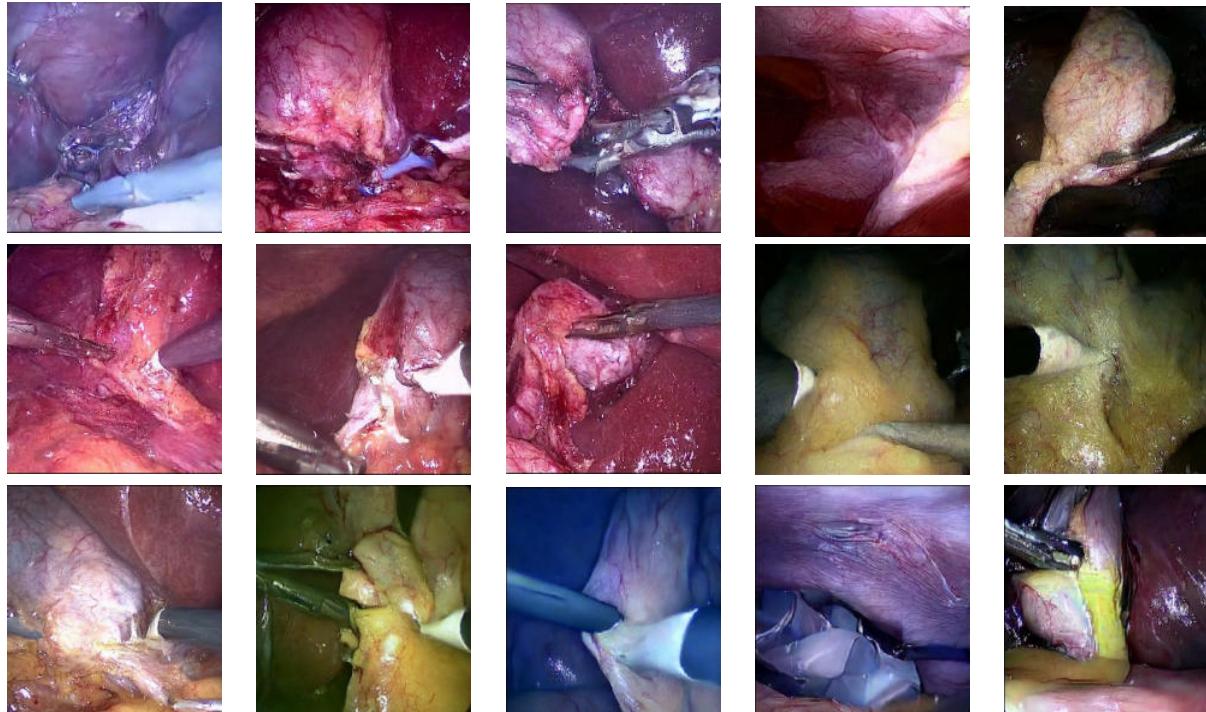
How can Generative AI enhance surgical training? The data basis is not given.



“Please, create an image from a laparoscopic cholecystectomy surgery, which displays a surgical grasper grasping a gallbladder”

Images are created with Dall-E of OpenAI

How can Generative AI enhance surgical training? Training and Adaption of Foundation Model.



[20]

Allmendinger, Hemmer, Queisner, Sauer, Kühl et al., (2023). Navigating the Synthetic Realm: Harnessing Diffusion-based Models for Laparoscopic Text-to-Image Generation [20]

How can Generative AI enhance surgical training? Training and Adaption of Foundation Model.

*“grasper grasp
gallbladder in
phase of
gallbladder
removal”*



Tool Movement



<https://anonymous.4open.science/w/laparoscopic-video-generation-D1C3/>

A perspective on Foundation Models

How should we adapt Foundation Models as a non-tech company?

Discussion



- 1 Introduction
- 2 Multimodality in Data
- 3 Generative Foundation Models
- 4 Research
- 5 Industry

Stable Diffusion

A glimpse into the battle for AI concentration

stability.ai

Stable Diffusion



Handelsblatt

Emad Mostaque kritisiert
„Machtkonzentration in der KI“

[21]

Wegen "zentralisierter KI": Stability AI verliert CEO

"Alles für eine dezentrale KI" – der bisherige CEO von Stability AI wendet sich gegen die Neuausrichtung des Unternehmens.



heise online

Bilder, die mit Stable Diffusion generiert wurden. (Bild: Stability AI)

[22]

<https://www.handelsblatt.com/technik/ki/kuenstliche-intelligenz-der-chef-von-einem-von-europas-wichtigsten-ki-unternehmen-tritt-ab/100027492.html> [21]

<https://www.heise.de/news/Wegen-zentralisierter-KI-Stability-AI-verliert-CEO-9664706.html> [22]

Who pays the bill for model training?

“No idea how we may one day generate revenue...”

GPT-4 Turbo

With 128k context, fresher knowledge and the broadest set of capabilities, GPT-4 Turbo is more powerful than GPT-4 and offered at a lower price.

[Learn about GPT-4 Turbo ↗](#)

Model	Input	Output
gpt-4-turbo-2024-04-09	US\$10.00 / 1M tokens	US\$30.00 / 1M tokens
Vision pricing calculator		
Set width	Set height	
512 <input type="button" value="px"/> by <input type="button" value="px"/>	=	US\$0.00255 <input type="button" value="i"/>
<input type="checkbox"/> Low resolution		
Price per 1K tokens (fixed)	US\$0.01	
512 × 512 tiles	1 × 1	
Total tiles	1	
Base tokens	85	
Tile tokens	170 × 1 = 170	
Total tokens	255	
Total price	US\$0.00255	

[23]



Technology

OpenAI hits \$2 bln revenue milestone - FT

By Reuters

February 9, 2024 11:52 AM GMT+1 - Updated 3 months ago



REUTERS [24]

<https://openai.com/api/pricing> (7/5/2024) [23]
<https://www.reuters.com/technology/openai-hits-2-bln-revenue-milestone-ft-2024-02-09/> [24]

What was the Gemini Controversy?

Google forced Gemini to temporarily stop users from creating pictures of people



<https://media.cnn.com/api/v1/images/stellar/prod/02222024-ai-generated-image-popegemini.jpg?c=original> [25]

What is NVIDIA up to?
Generating photorealistic simulated environments.



Isaac Sim

<https://developer.nvidia.com/isaac/sim> [26]

What is NVIDIA up to?

Foundation Models train themselves in simulated environments.



Isaac Sim



Eureka

<https://developer.nvidia.com/isaac/sim> [26]
Ma et al., (2023). Eureka: Human-Level Reward Design via Coding Large Language Models [27]

Why does Google want to be a part of it?
Foundation Models simulate realistic response to action.

Simulating long sequence of robot executions.

Step 1:



[28]

<https://universal-simulator.github.io/unisim/> [28]

What is NVIDIA up to? Foundation Models help to bring the simulation into reality.



Ma et al., (2024). DrEureka: Language Model Guided Sim-to-Real Transfer. [29]

Examples from industry

Foundation Models for visuals are already implemented as business asset.



[30]



[31]



[32]

Community

Engagement

Branding

<https://creatrealmagic.com> [30]

<https://aws.amazon.com/de/personalize/> [31]

<https://www.bmw.com/de/innovation/kreative-ai-bmw-8er-gran-coupe-kunstwerk-mit-kuenstlicher-intelligenz.html> [32]

A perspective on Foundation Models

Is there a concentration of power in AI - and is it bad?

Discussion